

Panning for Gold: Model-free Knockoffs for High-dimensional Controlled Variable Selection

Emmanuel Candès^{*1}, Yingying Fan², Lucas Janson¹, and Jinchi Lv²

¹Department of Statistics, Stanford University

²Data Sciences and Operations Department, Marshall School of Business, USC

Abstract

A common problem in modern statistical applications is to select, from a large set of candidates, a subset of variables which are important for determining an outcome of interest. For instance, the outcome may be disease status and the variables may be hundreds of thousands of single nucleotide polymorphisms on the genome. For data coming from low-dimensional ($n \geq p$) linear homoscedastic models, the knockoff procedure recently introduced by Barber and Candès solves the problem by performing variable selection while controlling the false discovery rate (FDR). The present paper extends the knockoff framework to arbitrary (and unknown) conditional models and any dimensions, including $n < p$, allowing it to solve a much broader array of problems. This extension requires the design matrix be random (independent and identically distributed rows) with a covariate distribution that is known, although we show our procedure to be robust to unknown/estimated distributions. To our knowledge, no other procedure solves the variable selection problem in such generality, but in the restricted settings where competitors exist, we demonstrate the superior power of knockoffs through simulations. Finally, we apply our procedure to data from a case-control study of Crohn’s disease in the United Kingdom, making twice as many discoveries as the original analysis of the same data.

Keywords. False discovery rate (FDR), knockoff filter, testing for conditional independence in nonlinear models, generalized linear models, logistic regression, Markov blanket, genome-wide association study (GWAS)

1 Introduction

1.1 Panning for gold

Certain diseases have a genetic basis, and an important biological problem is to find which genetic features (e.g., gene expressions or single nucleotide polymorphisms) are important for determining a given disease. In health care, researchers often want to know which electronic medical record entries determine future medical costs. Political scientists study which demographic or socioeconomic variables determine political opinions. Economists are similarly interested in which demographic/socioeconomic variables affect future income. Those in the technology industry seek out specific software characteristics they can change to increase user engagement. In the current data-driven science and engineering era, a list of such problems would go on and on. The common theme in all these instances is that we have a deluge of explanatory variables, often many more than the number of observations, knowing full well that the outcome we wish to understand better only actually depends on a small fraction of them. Therefore, a primary goal in modern “big data analysis” is to identify those important predictors in a sea of noise variables. Having said this, a reasonable question is why do we have so many covariates in the first place? The answer is twofold: first, because we can. To be sure, it may be fairly easy to measure thousands if not millions of attributes at the same time. For instance, it has become relatively inexpensive to genotype an individual, collecting hundreds of thousands of genetic variations at once. Second, even though we may believe that a trait or phenotype depends on a comparably small set of genetic variations, we have a priori no idea about which ones are relevant and therefore must include them all in our search for those nuggets of gold, so to speak. To further complicate matters, a common challenge in these big data problems, and a central focus of this paper, is that we often have little to no knowledge of how the outcome even depends on the few truly important variables.

To cast the ubiquitous (*model*) *selection* problem in statistical terms, call Y the random variable representing the outcome whose determining factors we are interested in, and X_1, \dots, X_p the set of p potential explanatory variables.

^{*}Author names are sorted alphabetically.

The object of study is the *conditional* distribution of the outcome Y given the covariates $X = (X_1, \dots, X_p)$, and we shall denote this conditional distribution function by $F_{Y|X}$. Ideally we would like to estimate $F_{Y|X}$, but in general, this is effectively impossible from a finite sample. For instance, even knowing that the conditional density depends upon 20 *known* covariates makes the problem impossible unless either the sample size n is astronomically large, and/or we are willing to impose a very restrictive model. However, in most problems $F_{Y|X}$ may realistically be assumed to depend on a small fraction of the p covariates; that is, the function $F_{Y|X}(y|x_1, \dots, x_p)$ only depends upon a small number of coordinates x_i (or is well approximated by such a lower-dimensional function). Although this assumption does not magically make the estimation of $F_{Y|X}$ easy, it does suggest consideration of the simpler problem: *which of the many variables does Y depend upon?* Often, finding a few of the important covariates—in other words, teasing out the relevant factors from those which are not—is already scientifically extremely useful and can be considered a first step in understanding the dependence between an outcome and some interesting variables; we regard this as a crucial problem in modern data science.

1.2 A peek at our contribution

This paper addresses the selection problem by considering a very general conditional model, where the response Y can depend in an arbitrary fashion on the covariates X_1, \dots, X_p . The only restriction we place on the model is that the observations $(X_{i1}, \dots, X_{ip}, Y_i)$ are independently and identically distributed (i.i.d.), which is often realistic in high-dimensional applications such as genetics, where subjects may be drawn randomly from some large population, or client behavioral modeling, where experiments on a service or user interface go out to a random subset of users. Therefore, the model is simply

$$(X_{i1}, \dots, X_{ip}, Y_i) \stackrel{i.i.d.}{\sim} F_{XY}, \quad i = 1, \dots, n \quad (1.1)$$

for some arbitrary $(p+1)$ -dimensional joint distribution F_{XY} . We will assume *no knowledge* of the conditional distribution of $Y | X_1, \dots, X_p$, but we do assume the joint distribution of the covariates is known, and we will denote it by F_X .

In Section 1.3 below we shall discuss the merits of this model but we would like to immediately remark on an important benefit: namely, one can pose a meaningful problem. To do this, observe that when we say that the conditional distribution of Y actually depends upon a (small) subset $\mathcal{S} \subset \{1, \dots, p\}$ of the variables X_1, \dots, X_p , which we would like to identify, we mean that we would like to find the “smallest” subset \mathcal{S} such that conditionally on $\{X_j\}_{j \in \mathcal{S}}$, Y is independent of all other variables. Another way to say this is that the other variables do not provide additional information about Y . A minimal set \mathcal{S} with this property is usually called a Markov blanket or Markov boundary for Y in the literature on graphical models (Pearl, 1988, Section 3.2.1). Under very mild conditions about the joint distribution F_{XY} , the Markov blanket is well defined and unique (see Section 2 for details) so that we have a cleanly stated selection problem.

In most problems of interest, even with the knowledge of F_X , it is beyond hope to recover the blanket \mathcal{S} with no error. Hence, we are naturally interested in procedures that control a Type I error; that is, we would like to find as many variables as possible while at the same time not having too many false positives. In this paper, we focus on controlling the false discovery rate (FDR) (Benjamini and Hochberg, 1995), which we can define here as follows: letting $\hat{\mathcal{S}}$ be the outcome of a selection procedure operating on the sampled data (we put a hat because $\hat{\mathcal{S}}$ is random), the FDR is

$$\text{FDR} := \mathbb{E}[\text{FDP}], \quad \text{FDP} = \frac{\#\{j : j \in \hat{\mathcal{S}} \setminus \mathcal{S}\}}{\#\{j : j \in \hat{\mathcal{S}}\}} \quad (1.2)$$

with the convention $0/0 = 0$. Procedures that control the FDR are interpretable, as they roughly bound what fraction of discoveries are false ones, and they can be quite powerful as well.

One achievement of this paper is to show that we can design rather powerful procedures that rigorously control the FDR (1.2) in finite samples. This holds no matter the unknown relationship between the explanatory variables X and the outcome Y . We achieve this by building upon the work of Barber and Candès (2015), who originally introduced the knockoff framework (throughout this paper, we will sometimes use “knockoffs” as shorthand for the knockoff framework or procedure). Their salient idea was to construct a set of so-called “knockoff” variables which were not (conditionally on the original variables) associated with the response, but whose structure mirrored that of the original covariates. These knockoff variables could then be used as controls for the real covariates, so that only real covariates which appeared to be considerably more associated with the response than their knockoff counterparts were selected. Their main result was achieving exact finite-sample FDR control, conditional on the observed covariates (so no random design assumption), in the homoscedastic Gaussian linear regression model when $n \geq 2p$, along with a nearly-exact extension to when $p \leq n < 2p$.

In this work, we apply the “knockoff philosophy” to the general random design model we have introduced, allowing us to extend Barber and Candès (2015) to arbitrary models and remove the low-dimensional constraint. We refer to the new approach as *model-free* (MF) knockoffs. In a nutshell:

- We propose a new knockoff construction amenable to the random design setting.
- As in Barber and Candès (2015) and further reviewed in Section 3, we shall use the knockoff variables as controls in such a way that we can tease apart important variables from noise while controlling the FDR. An appealing feature of this method is that the procedure completely bypasses the need for p-values.
- We apply the new procedure to real data from a case-control study of Crohn’s disease in the United Kingdom, please see Section 6. There, we show that the new knockoff method makes twice as many discoveries as the original analysis of the same data.

Before turning to the presentation of our method and results, we pause to discuss the merits and limitations of our model, the relationships between this work and others on selective inference, and the larger problem of high-dimensional statistical testing. The rest of this section is devoted to these points.

1.3 Relationship with the classical setup for inference

It may seem to the statistician that our model appears rather different from what she is used to. Our framework is, however, not as exotic as it looks.

Classical setup The usual setup for inference in conditional models is to assume a strong parametric model for the response conditional on the covariates, such as a homoscedastic linear model, but to assume as little as possible about, or even condition on, the covariates. We do the exact opposite by assuming we know *everything* about the covariate distribution but *nothing* about the conditional distribution $Y|X_1, \dots, X_p$. In practice, the parametric model in the classical approach is just an approximation, and does not need to hold exactly to produce useful inference. Analogously, we do not need to know the covariate distribution exactly for our method to be useful, as we will demonstrate in Sections 5.4 and 6.

When are our assumptions useful? We do not claim our assumptions will always be appropriate, but there are important cases when it is reasonable to think we know much more about the covariate distribution than about the conditional distribution of the response, including:

- When we in fact know exactly the covariate distribution because we control it, such as in gene knockout experiments (Cong et al., 2013; Peters et al., 2016), genetic crossing experiments (Haldane and Waddington, 1931), or sensitivity analysis of numerical models (Saltelli et al., 2008) (for example climate models). In some cases we may also essentially know the covariate distribution even when we do not control it, such as in admixture mapping (Tang et al., 2006).
- When we have a large amount of unsupervised data (covariate data without corresponding responses/labels) in addition to the n labeled observations. This is not uncommon in genetic or economic studies, where many other studies will exist that have collected the same covariate information but different response variables.
- When we simply have considerably more prior information about the covariates than about the response. Indeed, the point of many conditional modeling problems is to relate a poorly-understood response variable to a set of well-understood covariates. For instance, in genetic case-control studies, scientists seek to understand the causes of an extremely biologically-complex disease using many comparatively simple single nucleotide polymorphisms (SNPs) as covariates.

Payoff There are substantial payoffs to our framework. Particularly in high dimensions, previous inference results rely not only on a parametric model that is often linear and homoscedastic, but also on the sparsity or ultra-sparsity of the parameters of that model in order to achieve some asymptotic guarantee. In contrast, our framework can accommodate *any* model for both the response and the covariates, and our guarantees are exact in finite samples (non-asymptotic). In particular, our setup encompasses any regression, classification, or survival model, including any generalized linear model (GLM), and allows for arbitrary nonlinearities and heteroscedasticity, such as are found in many machine learning applications.

1.4 Relationship with work on inference after selection

There is a line of work on inference after selection, or post-selection inference, for high-dimensional regression, the goal of which is to first perform selection to make the problem low-dimensional, and then produce p-values that are valid *conditional* on the selection step (Berk et al., 2013; Lockhart et al., 2014; Lee et al., 2016). These works differ from ours in a number of ways so that we largely see them as complementary activities.

- First, our focus is on selecting the right variables, whereas the goal of this line of work is to adjust inference after some selection has taken place. In more detail, these works presuppose a selection procedure has been chosen (for reasons that may have nothing to do with controlling Type I error) and then compute p-values for the selected variables, taking into account the selection step. In contrast, MF knockoffs is by itself a selection procedure that controls Type I error.
- Second, inference after selection relies heavily on parametric assumptions about the conditional distribution, namely that

$$Y|X_1, \dots, X_p \sim \mathcal{N}(\mu(X_1, \dots, X_p), \sigma^2),$$

making it unclear how to extend it to the more general setting of the present paper.

- The third difference stems from their objects of inference. In the selection step, a subset of size $m \leq n$ of the original p covariates is selected, say X_{j_1}, \dots, X_{j_m} , and the objects of inference are the coefficients of the X_{j_k} 's in the projection of μ onto the linear subspace spanned by the $n \times m$ matrix of observed values of these X_{j_k} 's. That is, the k th null hypothesis is that the aforementioned coefficient on X_{j_k} is zero—note that whether or not inference on the j th variable is produced at all, and if it is, the object of that inference, both depend on the initial selection step. In contrast, if MF knockoffs were restricted to the homoscedastic Gaussian model above, the j th null hypothesis would be that μ does not depend on X_j , and there would be p such null hypotheses, one for each of the original variables.

1.5 Obstacles to obtaining p-values

Our procedure does not follow the canonical approach to FDR control and multiple testing in general. The canonical approach is to plug p-values into the BHq procedure, which controls the FDR under p-value independence and certain forms of dependence (Benjamini and Hochberg, 1995; Benjamini and Yekutieli, 2001). Although these works have seeded a wealth of methodological innovations over the past two decades (Benjamini, 2010), all these procedures act on a set of valid p-values (or equivalent statistics), which they assume can be computed.¹ The requirement of having valid p-values is quite constraining for general conditional modeling problems.

1.5.1 Regression p-value approximations

In low-dimensional ($n \geq p$) homoscedastic Gaussian linear regression, p-values can be computed exactly even if the error variance is unknown, although the p-values will not in general have any simple dependence properties like independence or positive regression dependency on a subset (PRDS). Already for just the slightly broader class of low-dimensional GLMs, one must resort to asymptotic p-values derived from maximum-likelihood theory, which we will show in Section 5.1 can be far from valid in practice. In high-dimensional ($n < p$) GLMs, it is not clear how to get p-values at all. Although some work (see for example van de Geer et al. (2014)) exists on computing asymptotic p-values under strong sparsity assumptions, like their low-dimensional maximum-likelihood counterparts, these methods suffer from highly non-uniform null p-values in many finite-sample problems. For binary covariates, the causal inference literature uses matching and propensity scores for approximately valid inference, but extending these methods to high dimensions is still a topic of current research, requiring similar assumptions and asymptotic approximations to the aforementioned high-dimensional GLM literature (Athey et al., 2016). Moving beyond generalized linear models to the nonparametric setting, there exist measures of feature importance, but no p-values.²

¹One notable exception is the SLOPE procedure (Bogdan et al., 2015), which acts like a regression analogue of BHq without ever computing p-values, and provably controls the FDR in homoscedastic linear regression when the design matrix has orthogonal columns (necessitating, importantly, that $n \geq p$) and empirically, otherwise, whenever the signal obeys sparsity constraints.

²In their online description of random forests (<http://www.math.usu.edu/~adele/forests/>), Leo Breiman and Adele Cutler propose a way to obtain a “z-score” for each variable, but without any theoretical distributional justification, and Strobl and Zeileis (2008) find “that the suggested test is not appropriate for statements of significance.”

1.5.2 Marginal testing

Faced with the inability to compute p-values for hypothesis tests of conditional independence, one solution is to use *marginal* p-values, i.e., p-values for testing *unconditional* (or marginal) independence between Y and X_j . This simplifies the problem considerably, and many options exist for obtaining valid p-values for such a test. However, marginal p-values are in general *invalid* for testing conditional independence, and replacing tests of conditional independence with tests of unconditional independence is often undesirable. Indeed when $p \ll n$, so that classical (e.g., maximum-likelihood) inference techniques for regression give valid p-values for parametric tests of conditional independence, it would be very unusual to resort to marginal testing to select important covariates, and we cannot think of a textbook that takes this route. Furthermore, the class of conditional test statistics is far richer than that of marginal ones, and includes the most powerful statistical inference and prediction methodology available. For example, in compressed sensing, the signal recovery guarantees for state-of-the-art ℓ_1 -based (joint) algorithms are stronger than any guarantees possible with marginal methods. To constrain oneself to marginal testing is to completely ignore the vast modern literature on sparse regression that, while lacking finite-sample Type I error control, has had tremendous success establishing other useful inferential guarantees such as model selection consistency under high-dimensional asymptotics in both parametric (e.g., Lasso (Zhao and Yu, 2006; Candès and Plan, 2009)) and nonparametric (e.g., random forests (Wager and Athey, 2016)) settings. Realizing this, the statistical genetics community has worked on a number of multivariate approaches to improve power in genome-wide association studies using both penalized (Wu et al., 2009; He and Lin, 2011) and Bayesian regression (Guan and Stephens, 2011; Li et al., 2011), but both approaches still suffer from a lack of Type I error control (without making strong assumptions on parameter priors). We will see that the MF knockoff procedure is able to leverage the power of any of these techniques while adding rigorous finite-sample Type I error control.

Some specific drawbacks of using marginal p-values are:

1. **Power loss** Even when the covariates are independent, so that the hypotheses of conditional and unconditional independence coincide, p-values resulting from marginal testing procedures may be less powerful than those from conditional testing procedures. This phenomenon has been reported previously, for example in statistical genetics by Hoggart et al. (2008) and many others. Intuitively, this is because a joint model in X_1, \dots, X_p for Y will have less residual variance than a marginal model in just X_j . There are exceptions, for instance if there is only one important variable, then its marginal model is the correct joint model and a conditional test will be less powerful due to the uncertainty in how the other variables are included in the joint model. But in general, marginal testing becomes increasingly underpowered relative to conditional testing as the *absolute* number of important covariates increases (Frommlet et al., 2012), suggesting particular advantage for conditional testing in modern applications with complex high-dimensional models.

There are also cases in which important variables are in fact fully marginally independent of the response. As a toy example, if $X_1, X_2 \stackrel{\text{iid}}{\sim} \text{Bernoulli}(0.5)$ and $Y = \mathbb{1}_{\{X_1 + X_2 = 1\}}$, then Y is marginally independent of each of X_1 and X_2 , even though together they determine Y perfectly. X_1 and X_2 are certainly *conditionally* dependent on Y , however, so a conditional test can have power to discover them.

2. **Interpretability** When the covariates are not independent, marginal and conditional independence do not correspond, and we end up asking the wrong question. For example, in a model with a few important covariates which cause the outcome to change, and many unimportant covariates which have no influence on the outcome but are correlated with the causal covariates, marginal testing will treat such unimportant covariates as important. Thus, because marginal testing is testing the wrong hypotheses, there will be many “discoveries” which have no influence on the outcome.

The argument is often made, especially in genetics, that although discovering a non-causal covariate just because it was correlated with a causal one is technically incorrect, it can still be useful as it suggests that there is a causal covariate correlated with the discovered one. While this is indeed useful, especially in genetics where correlated SNPs tend to be very close to one another on the genome, this comes at a price since it significantly alters the meaning of the FDR. Indeed, if we adopt this viewpoint, the unit of inference is no longer a SNP but, rather, a region on the genome, yet FDR is still being tested at the SNP level. A consequence of this mismatch is that the result of the analysis may be completely misleading, as beautifully argued in Pacifico et al. (2004); Benjamini and Heller (2007); Siegmund et al. (2011), see also Chouldechova (2014) and Brzyski et al. (2016) for later references.

3. **Dependent p-values** Marginal p-values will, in general, have quite complicated joint dependence, so that BHq does not control FDR exactly. Although procedures for controlling FDR under arbitrary dependence exist, their increased generality tends to make them considerably more conservative than BHq. In practice, however,

the FDR of BHq applied to dependent p-values is usually below its nominal level, but the problem is that it can have highly *variable* FDP. Recall that FDR is the expectation of FDP, the latter being the random quantity we actually care about but cannot control directly. Therefore, FDR control is only useful if the realized FDP is relatively concentrated around its expectation, and it is well-established (Efron, 2010, Chapter 4) that under correlations BHq can produce highly skewed FDP distributions. In such cases, with large probability, $FDP = 0$ perhaps because no discoveries are made, and when discoveries are made, the FDP may be much higher than the nominal FDR, making it a misleading error bound.

1.6 Getting valid p-values via conditional randomization testing

If we insist on obtaining p-values for each X_j , there is in fact a simple method when the covariate distribution is assumed known, as it is in this paper. This method is similar in spirit to both propensity scoring (where the distribution of a binary X_j conditional on the other variables is often estimated) and randomization/permutation tests (where X_j is either the only covariate or fully independent of the other explanatory variables), exists. Explicitly, a conditional randomization test for the j th variable proceeds by first computing some feature importance statistic T_j for the j th variable. Then the null distribution of T_j can be computed through simulation by independently sampling X_j^* 's from the *conditional* distribution of X_j given the others (derived from the known F_X) and recomputing the same statistic T_j^* with each new X_j^* in place of X_j , see Section 4 for details. Despite its simplicity, we have not seen this test proposed previously in the literature, although it nearly matches the usual randomization test when the covariates are independent of one another.

1.7 Outline of the paper

The remainder of the paper is structured as follows:

- Section 2 frames the controlled selection problem in rigorous mathematical terms.
- Section 3 introduces the MF knockoff procedure, examines its relationship with the earlier proposal of Barber and Candès (2015), proposes knockoff constructions and feature statistics, and establishes FDR control.
- Section 4 introduces the conditional randomization test.
- Section 5 demonstrates through simulations that the MF knockoff procedure controls the FDR in a number of settings where no other procedure does, and that when competitors exist, knockoffs is more powerful.
- Section 6 applies our procedure to a case-control study of Crohn's disease in the United Kingdom. Using the design matrix from the same data set, we also show that knockoffs' FDR control is surprisingly robust to covariate distribution estimation error, making it useful even when the covariate distribution is not known exactly.
- Section 7 concludes the paper with extensions and potential lines of future research.

2 Problem statement

To state the controlled variable selection problem carefully, suppose we have n i.i.d. samples from a population, each of the form (X, Y) , where $X = (X_1, \dots, X_p) \in \mathbb{R}^p$ and $Y \in \mathbb{R}$. If the conditional distribution of Y actually depends upon a smaller subset of these variables, we would like to classify each variable as relevant or not depending on whether it belongs to this subset or not. Mathematically speaking, we are looking for the Markov blanket \mathcal{S} , i.e. the “smallest” subset \mathcal{S} such that conditionally on $\{X_j\}_{j \in \mathcal{S}}$, Y is independent of all other variables. For almost all joint distributions of (X, Y) , there exists a unique Markov blanket but there are pathological cases where it does not. An example is this: suppose that X_1 and X_2 are independent Gaussian variables and that $X_3 = X_1 - X_2$. Further assume that the distribution of Y depends upon the vector X only through $X_1 + X_2$, e.g., $Y|X \sim \mathcal{N}(X_1 + X_2, 1)$. Then the set of relevant variables—or equivalently, the Markov blanket—is ill defined since we can say that the likelihood of Y depends upon X through either (X_1, X_2) , (X_1, X_3) , or (X_2, X_3) , all these subsets being equally good. In order to define a unique set of relevant variables, we shall work with the notion of conditional *pairwise* independence.

Definition 2.1. A variable X_j is said to be “null” if and only if Y is independent of X_j conditionally on the other variables $X_{-j} = \{X_1, \dots, X_p\} \setminus \{X_j\}$. The subset of null variables is denoted by $\mathcal{H}_0 \subset \{1, \dots, p\}$ and we call a variable X_j “nonnull” or relevant if $j \notin \mathcal{H}_0$.

From now on, *our goal is to discover as many relevant (conditionally dependent) variables as possible while keeping the FDR under control.*³ Formally, for a selection rule that selects a subset $\hat{\mathcal{S}}$ of the covariates,

$$\text{FDR} := \mathbb{E} \left[\frac{|\hat{\mathcal{S}} \cap \mathcal{H}_0|}{|\hat{\mathcal{S}}|} \right]. \quad (2.1)$$

In the example above, because of the perfect functional relationship $X_3 = X_2 - X_1$, all three variables X_1, X_2, X_3 would be classified as nulls. Imagine, however, breaking this relationship by adding a bit of noise, e.g., $X_3 = X_2 - X_1 + Z$, where Z is Gaussian noise (independent of X_1 and X_2) however small. Then according to our definition, X_1 and X_2 are both nonnull while X_3 is null—and everything makes sense. Having said this, we should not let ourselves be distracted by such subtleties. In the literature on graphical models there, in fact, exist weak regularity conditions that guarantee that the (unique) set of relevant variables defined by pairwise conditional independence, exactly coincides with the Markov blanket so that there is no ambiguity. In this field, researchers typically assume these weak regularity conditions hold (examples would include the local and global Markov properties), and proceed from there. For example, the textbook Edwards (2000) describes these properties on page 8 as holding “under quite general conditions” and then assumes them for the rest of the book.

Our definition is very natural to anyone working with parametric GLMs. In a GLM, the response Y has a probability distribution taken from an exponential family, which depends upon the covariates only through the linear combination $\eta = \beta_1 X_1 + \dots + \beta_p X_p$. The relationship between Y and X is specified via a link function g such that $\mathbb{E}(Y|X) = g^{-1}(\eta)$. In such models and under broad conditions, $Y \perp\!\!\!\perp X_j | X_{-j}$ if and only if $\beta_j = 0$. In this context, testing the hypothesis that X_j is a null variable is the same as testing $H_j : \beta_j = 0$.

Proposition 2.2. *Take a family of random variables X_1, \dots, X_p such that one cannot perfectly predict any one of them from knowledge of the others. If the likelihood of Y follows a GLM, then $Y \perp\!\!\!\perp X_j | X_{-j}$ if and only if $\beta_j = 0$. Hence, \mathcal{H}_0 from Definition 2.1 is exactly the set $\{j : \beta_j = 0\}$.*

Proof. We prove this in the case of the logistic regression model as the general case is similar. Here, the conditional distribution of Y is Bernoulli with

$$\mathbb{E}(Y|X) = \mathbb{P}(Y = 1|X) = \frac{e^\eta}{1 + e^\eta} = g^{-1}(\eta), \quad \eta = \beta_1 X_1 + \dots + \beta_p X_p,$$

and please note that the assumption about the covariates implies that the model is identifiable. Now assume first that $\beta_j = 0$. Then

$$p_{Y, X_j | X_{-j}}(y, x_j | x_{-j}) = p_{Y | X_j, X_{-j}}(y | x_j, x_{-j}) p_{X_j | X_{-j}}(x_j | x_{-j}) \quad (2.2)$$

and since the first factor in the right-hand side does not depend on X_j , we see that the conditional probability distribution function factorizes. This implies conditional independence. In the other direction, assume that Y and X_j are conditionally independent. Then the likelihood function

$$\frac{\exp(Y(\beta_1 X_1 + \dots + \beta_p X_p))}{1 + \exp(\beta_1 X_1 + \dots + \beta_p X_p)}$$

must, conditionally on X_{-j} , factorize into a function of Y times a function of X_j . A consequence of this is that conditionally on X_{-j} , the odds ratio must not depend on X_j (it must be constant). However, this ratio is equal to $\exp(\beta_j X_j)$ and is constant only if $\beta_j = 0$ since, by assumption, X_j is not determined by X_{-j} . \square

The assumption regarding the covariates is needed. Indeed, suppose $X_1 \sim \mathcal{N}(0, 1)$, $X_2 = 1\{X_1 > 0\}$, and Y follows a logistic model as above with $\eta = X_1 + X_2$. Then $Y \perp\!\!\!\perp X_2 | X_1$ even though $\beta_2 = 1$. In this example, the conditional distribution of Y depends on (X_1, X_2) only through X_1 . Therefore, for the purpose of identifying important variables (recall our task is to find important variables and not to learn exactly how the likelihood function depends upon these variables), we would like to find X_1 and actually do not care about X_2 since it provides no new information.

³Using the methods of Janson and Su (2016), other error rates such as the k -familywise error rate can also be controlled using MF knockoffs, but we focus on FDR for this paper.

3 Methodology

3.1 Model-free knockoffs

3.1.1 Definition

Definition 3.1. *Model-free knockoffs for the family of random variables $X = (X_1, \dots, X_p)$ are a new family of random variables $\tilde{X} = (\tilde{X}_1, \dots, \tilde{X}_p)$ constructed with the following two properties: (1) for any subset $S \subset \{1, \dots, p\}$,⁴*

$$(X, \tilde{X})_{\text{swap}(S)} \stackrel{d}{=} (X, \tilde{X}); \quad (3.1)$$

(2) $\tilde{X} \perp\!\!\!\perp Y \mid X$ if there is a response Y . (2) is guaranteed if \tilde{X} is constructed without looking at Y .

Above, the vector $(X, \tilde{X})_{\text{swap}(S)}$ is obtained from (X, \tilde{X}) by swapping the entries X_j and \tilde{X}_j for each $j \in S$; for example, with $p = 3$ and $S = \{2, 3\}$,

$$(X_1, X_2, X_3, \tilde{X}_1, \tilde{X}_2, \tilde{X}_3)_{\text{swap}(\{2,3\})} = (X_1, \tilde{X}_2, \tilde{X}_3, \tilde{X}_1, X_2, X_3).$$

We see from (3.1) that original and knockoff variables are pairwise exchangeable: taking any subset of variables and swapping them with their knockoffs leaves the joint distribution invariant. To give an example of MF knockoffs, suppose that $X \sim \mathcal{N}(0, \Sigma)$. Then a joint distribution obeying (3.1) is this:

$$(X, \tilde{X}) \sim \mathcal{N}(0, \mathbf{G}), \quad \text{where} \quad \mathbf{G} = \begin{bmatrix} \Sigma & \Sigma - \text{diag}\{s\} \\ \Sigma - \text{diag}\{s\} & \Sigma \end{bmatrix}; \quad (3.2)$$

here, $\text{diag}\{s\}$ is any diagonal matrix selected in such a way that the joint covariance matrix \mathbf{G} is positive semidefinite. Indeed, the distribution obtained by swapping variables with their knockoffs is Gaussian with a covariance given by $\mathbf{P}\mathbf{G}\mathbf{P}$, where \mathbf{P} is the permutation matrix encoding the swap. Since $\mathbf{P}\mathbf{G}\mathbf{P} = \mathbf{G}$ for any swapping operation, the distribution is invariant.

We will soon be interested in the problem of constructing knockoff variables, having observed X . In the above example, a possibility is to sample the knockoff vector \tilde{X} from the conditional distribution

$$\tilde{X} \mid X \stackrel{d}{=} \mathcal{N}(\mu, \mathbf{V}),$$

where μ and \mathbf{V} are given by classical regression formulas, namely,

$$\begin{aligned} \mu &= X - X \text{diag}\{s\} \Sigma^{-1}, \\ \mathbf{V} &= 2 \text{diag}\{s\} - \text{diag}\{s\} \Sigma^{-1} \text{diag}\{s\}. \end{aligned}$$

There are, of course, many other ways of constructing knockoff variables, and for the time being, we prefer to postpone the discussion of more general constructions.

In the setting of the paper, we are given i.i.d. pairs $(X_{i1}, \dots, X_{ip}, Y_i) \in \mathbb{R}^p \times \mathbb{R}$ of covariates and responses, which we can assemble in a data matrix \mathbf{X} and a data vector y in such a way that the i th row of \mathbf{X} is (X_{i1}, \dots, X_{ip}) and the i th entry of y is Y_i . Then the MF knockoff matrix $\tilde{\mathbf{X}}$ is constructed in such a way that for each observation label i , $(\tilde{X}_{i1}, \dots, \tilde{X}_{ip})$ is a knockoff for (X_{i1}, \dots, X_{ip}) as explained above; that is to say, the joint vector $(X_{i1}, \dots, X_{ip}, \tilde{X}_{i1}, \dots, \tilde{X}_{ip})$ obeys the pairwise exchangeability property (3.1).

3.1.2 Relationship with the knockoffs of Barber and Candès

We can immediately see the key difference with the earlier framework of Barber and Candès (2015). In their work, the design matrix is viewed as being fixed and setting $\Sigma = \mathbf{X}^\top \mathbf{X}$, knockoff variables are constructed as to obey

$$[\mathbf{X}, \tilde{\mathbf{X}}]^\top [\mathbf{X}, \tilde{\mathbf{X}}] = \mathbf{G}, \quad \mathbf{G} = \begin{bmatrix} \Sigma & \Sigma - \text{diag}\{s\} \\ \Sigma - \text{diag}\{s\} & \Sigma \end{bmatrix}. \quad (3.3)$$

Imagine the columns of \mathbf{X} are centered, i.e., have vanishing means, so that we can think of $\mathbf{X}^\top \mathbf{X}$ as a sample covariance matrix. Then the construction of Barber and Candès is asking that the *sample* covariance matrix of the joint set of variables obeys the exchangeability property—i.e., swapping rows and columns leaves the covariance

⁴ $\stackrel{d}{=}$ denotes equality in distribution, and the definition of the swapping operation is given just below.

invariant—whereas in this paper, it is the *population* covariance that must be invariant. In particular, the MF knockoffs will be far from obeying the relationship $\mathbf{X}^\top \mathbf{X} = \tilde{\mathbf{X}}^\top \tilde{\mathbf{X}}$ required in (3.3). To drive this point home, assume $X \sim \mathcal{N}(0, \mathbf{I})$. Then we can choose $\tilde{X} \sim \mathcal{N}(0, \mathbf{I})$ independently from X , in which case $\mathbf{X}^\top \mathbf{X}/n$ and $\tilde{\mathbf{X}}^\top \tilde{\mathbf{X}}/n$ are two independent Wishart variables. An important consequence of this is that in the MF approach, the sample correlation between the j th columns of \mathbf{X} and $\tilde{\mathbf{X}}$ will typically be far smaller than that in the original framework of Barber and Candès. For example, take $n = 3000$ and $p = 1000$ and assume the equicorrelated construction from Barber and Candès (2015). Then the sample correlation between any variable \mathbf{X}_j and its knockoff $\tilde{\mathbf{X}}_j$ will be about 0.65 while in the random case, the average magnitude of the correlation is about 0.015. We shall see later how this explains the gain in power our new method brings about.

3.1.3 Exchangeability of null covariates and their knockoffs

A crucial property of MF knockoffs is that we can swap null covariates with their knockoffs without changing the joint distribution of the original covariates X , their knockoffs \tilde{X} , and the response Y . From now on, $X_{i:j}$ for $i \leq j$ is a shorthand for (X_i, \dots, X_j) .

Lemma 3.2. *Take any subset $S \subset \mathcal{H}_0$ of nulls. Then*

$$([\mathbf{X}, \tilde{\mathbf{X}}], y) \stackrel{d}{=} ([\mathbf{X}, \tilde{\mathbf{X}}]_{\text{swap}(S)}, y).$$

Proof. Assume without loss of generality that $S = \{1, 2, \dots, m\}$. By row independence, it suffices to show that $((X, \tilde{X}), Y) \stackrel{d}{=} ((X, \tilde{X})_{\text{swap}(S)}, Y)$, where X (resp. Y) is a row of \mathbf{X} (resp. y). Furthermore, since $(X, \tilde{X}) \stackrel{d}{=} (X, \tilde{X})_{\text{swap}(S)}$, we only need to establish that

$$Y \mid (X, \tilde{X})_{\text{swap}(S)} \stackrel{d}{=} Y \mid (X, \tilde{X}). \quad (3.4)$$

Letting $p_{Y|X}(y|x)$ be the conditional distribution of Y , observe that

$$\begin{aligned} p_{Y|(X, \tilde{X})_{\text{swap}(S)}}(y|(x, \tilde{x})) &= p_{Y|(X, \tilde{X})}(y|(x, \tilde{x})_{\text{swap}(S)}) \\ &= p_{Y|X}(y|x'), \end{aligned}$$

where $x'_i = \tilde{x}_i$ if $i \in S$ and $x'_i = x_i$ otherwise. The second equality above comes from the fact that Y is conditionally independent of \tilde{X} by property (2) in the definition of MF knockoffs. Next, since Y and X_1 are independent conditional on $X_{2:p}$, we have

$$\begin{aligned} p_{Y|X_{1:p}}(y|\tilde{x}_1, x'_{2:p}) &= p_{Y|X_{2:p}}(y|x'_{2:p}) \\ &= p_{Y|X_{1:p}}(y|x_1, x'_{2:p}). \end{aligned}$$

This shows that

$$Y \mid (X, \tilde{X})_{\text{swap}(S)} \stackrel{d}{=} Y \mid (X, \tilde{X})_{\text{swap}(S \setminus \{1\})}.$$

We can repeat this argument with the second variable, the third, and so on until S is empty. This proves (3.4). \square

3.2 Feature statistics

In order to find the relevant variables, we now compute statistics W_j for each $j \in \{1, \dots, p\}$, a large positive value of W_j providing evidence against the hypothesis that X_j is null. This statistic depends on the response and the original variables but also on the knockoffs; that is,

$$W_j = w_j([\mathbf{X}, \tilde{\mathbf{X}}], y)$$

for some function w_j . As in Barber and Candès (2015), we impose a *flip-sign property*, which says that swapping the j th variable with its knockoff has the effect of changing the sign of W_j . Formally, if $[\mathbf{X}, \tilde{\mathbf{X}}]_{\text{swap}(S)}$ is the matrix obtained by swapping columns in S ,

$$w_j([\mathbf{X}, \tilde{\mathbf{X}}]_{\text{swap}(S)}, y) = \begin{cases} w_j([\mathbf{X}, \tilde{\mathbf{X}}], y), & j \notin S, \\ -w_j([\mathbf{X}, \tilde{\mathbf{X}}], y), & j \in S. \end{cases} \quad (3.5)$$

In contrast to the aforementioned work, we do not require the sufficiency property that w_j depend on \mathbf{X} , $\tilde{\mathbf{X}}$, and y only through $[\mathbf{X}, \tilde{\mathbf{X}}]^\top [\mathbf{X}, \tilde{\mathbf{X}}]$ and $[\mathbf{X}, \tilde{\mathbf{X}}]^\top y$.

At this point, it may help the reader unfamiliar with the knockoff framework to think about knockoff statistics $W = (W_1, \dots, W_p)$ in two steps: first, consider a statistic T for each original and knockoff variable,

$$T \triangleq (Z, \tilde{Z}) = (Z_1, \dots, Z_p, \tilde{Z}_1, \dots, \tilde{Z}_p) = t([\mathbf{X}, \tilde{\mathbf{X}}], y),$$

with the idea that Z_j (resp. \tilde{Z}_j) measures the importance of X_j (resp. \tilde{X}_j). Assume the natural property that switching a variable with its knockoff simply switches the components of T in the same way, namely, for each $S \subset \{1, \dots, p\}$,

$$(Z, \tilde{Z})_{\text{swap}(S)} = t([\mathbf{X}, \tilde{\mathbf{X}}]_{\text{swap}(S)}, y). \quad (3.6)$$

Then one can construct a W_j obeying the flip-sign condition (3.5) by setting

$$W_j = f_j(Z_j, \tilde{Z}_j),$$

where f_j is any anti-symmetric function.⁵ (Conversely, any statistic W_j verifying the flip sign condition can be constructed in this fashion.) Adopting this approach, consider a regression problem and run the Lasso on the original design augmented with knockoffs,

$$\min_{b \in \mathbb{R}^{2p}} \frac{1}{2} \|y - [\mathbf{X}, \tilde{\mathbf{X}}]b\|_2^2 + \lambda \|b\|_1 \quad (3.7)$$

and denote the solution by $\hat{b}(\lambda)$ (the first p components are the coefficients of the original variables and the last p are for the knockoffs). Then the *Lasso coefficient-difference* (LCD) statistic sets $Z_j = |\hat{b}_j(\lambda)|$, $\tilde{Z}_j = |\hat{b}_{j+p}(\lambda)|$, and

$$W_j = Z_j - \tilde{Z}_j = |\hat{b}_j(\lambda)| - |\hat{b}_{j+p}(\lambda)|. \quad (3.8)$$

A large positive value of W_j provides some evidence that the distribution of Y depends upon X_j , whereas under the null W_j has a symmetric distribution and, therefore, is equally likely to take on positive and negative values, as we shall see next. Before moving on, however, please carefully observe that the value of λ in (3.8) does not need to be fixed in advance, and can be computed from y and $[\mathbf{X}, \tilde{\mathbf{X}}]$ in any data-dependent fashion as long as permuting the columns of \mathbf{X} does not change its value; for instance, it can be selected by cross-validation.

Lemma 3.3. *Conditional on $|W|$, the signs of the null W_j 's, $j \in \mathcal{H}_0$, are i.i.d. coin flips.*

Proof. Let $\epsilon = (\epsilon_1, \dots, \epsilon_p)$ be a sequence of independent random variables such that $\epsilon_j = \pm 1$ with probability 1/2 if $j \in \mathcal{H}_0$, and $\epsilon_j = 1$ otherwise. To prove the claim, it suffices to establish that

$$W \stackrel{d}{=} \epsilon \odot W, \quad (3.9)$$

where \odot denotes pointwise multiplication, i.e. $\epsilon \odot W = (\epsilon_1 W_1, \dots, \epsilon_p W_p)$. Now, take ϵ as above and put $S = \{j : \epsilon_j = -1\} \subset \mathcal{H}_0$. Consider swapping variables in S :

$$W_{\text{swap}(S)} \triangleq w([\mathbf{X}, \tilde{\mathbf{X}}]_{\text{swap}(S)}, y).$$

On the one hand, it follows from the flip-sign property that $W_{\text{swap}(S)} = \epsilon \odot W$. On the other hand, Lemma 3.2 implies that $W_{\text{swap}(S)} \stackrel{d}{=} W$ since $S \subset \mathcal{H}_0$. These last two properties give (3.9). \square

3.3 FDR control

From now on, our methodology follows that of Barber and Candès (2015) and we simply rehearse the main ingredients while referring to their paper for additional insights. It follows from Lemma 3.3 that the null statistics W_j are symmetric and that for any fixed threshold $t > 0$,

$$\#\{j : W_j \leq -t\} \geq \#\{\text{null } j : W_j \leq -t\} \stackrel{d}{=} \#\{\text{null } j : W_j \geq t\}.$$

Imagine then selecting those variables such that W_j is sufficiently large, e.g., $W_j \geq t$, then the false discovery proportion (FDP)

$$\text{FDP}(t) = \frac{\#\{\text{null } j : W_j \geq t\}}{\#\{j : W_j \geq t\}} \quad (3.10)$$

⁵An anti-symmetric function f is such that $f(v, u) = -f(u, v)$.

can be estimated via the statistic

$$\widehat{\text{FDP}}(t) = \frac{\#\{j : W_j \leq -t\}}{\#\{j : W_j \geq t\}}$$

since the numerator is an upward-biased estimate of the unknown numerator in (3.10). The idea of the knockoff procedure is to choose a data-dependent threshold as liberal as possible while having an estimate of the FDP under control. The theorem below shows that estimates of the FDR process can be inverted to give tight FDR control.

Theorem 3.4. *Choose a threshold $\tau > 0$ by setting⁶*

$$\tau = \min \left\{ t > 0 : \frac{\#\{j : W_j \leq -t\}}{\#\{j : W_j \geq t\}} \leq q \right\} \quad (\mathbf{Knockoffs}), \quad (3.11)$$

where q is the target FDR level (or $\tau = +\infty$ if the set above is empty). Then the procedure selecting the variables

$$\hat{S} = \{j : W_j \geq \tau\}$$

controls the modified FDR defined as

$$\text{mFDR} = \mathbb{E} \left[\frac{|\{j \in \hat{S} \cap \mathcal{H}_0\}|}{|\hat{S}| + 1/q} \right] \leq q.$$

The slightly more conservative procedure, given by incrementing the number of negatives by one,

$$\tau_+ = \min \left\{ t > 0 : \frac{1 + \#\{j : W_j \leq -t\}}{\#\{j : W_j \geq t\}} \leq q \right\} \quad (\mathbf{Knockoffs}+) \quad (3.12)$$

and setting $\hat{S} = \{j : W_j \geq \tau_+\}$, controls the usual FDR,

$$\mathbb{E} \left[\frac{|\{j \in \hat{S} \cap \mathcal{H}_0\}|}{|\hat{S}| \vee 1} \right] \leq q.$$

These results are non-asymptotic and hold no matter the dependence between the response and the covariates.

The proof is the same as that of Theorems 1 and 2 in Barber and Candès (2015)—and, therefore, omitted—since all we need is that the null statistics have signs distributed as i.i.d. coin flips. Note that Theorem 3.4 only tells one side of the story: Type I error control; the other very important side is power, which leads us to spend most of the remainder of the paper considering how best to construct knockoff variables and statistics.

3.4 Constructing model-free knockoffs

3.4.1 Exact constructions

We have seen in Section 3.1.1 one way of constructing MF knockoffs in the case where the covariates are Gaussian. How should we proceed for non-Gaussian data? In this regard, the characterization below may be useful.

Proposition 3.5. *The random variables $(\tilde{X}_1, \dots, \tilde{X}_p)$ are model-free knockoffs for (X_1, \dots, X_p) if and only if for any $j \in \{1, \dots, p\}$, the pair (X_j, \tilde{X}_j) is exchangeable conditional on all the other variables and their knockoffs (and, of course, $\tilde{X} \perp\!\!\!\perp Y | X$).*

The proof consists of simple manipulations of the definition and is, therefore, omitted. Our problem can thus also be posed as constructing pairs that are conditionally exchangeable. If the components of the vector X are independent, then any independent copy of X would work; that is, any vector \tilde{X} independently sampled from the same joint distribution as X would work. With dependent coordinates, we may proceed as follows:

Algorithm 1 Sequential Conditional Independent Pairs.

```

j = 1  while j ≤ p do
    | Sample  $\tilde{X}_j$  from  $\mathcal{L}(X_j | X_{-j}, \tilde{X}_{1:j-1})$ 
    | j = j + 1
end
```

⁶When we write $\min\{t > 0 : \dots\}$, we abuse notation since we actually mean $\min\{t \in \mathcal{W}_+ : \dots\}$, where $\mathcal{W}_+ = \{|W_j| : |W_j| > 0\}$.

Above, $\mathcal{L}(X_j | X_{-j}, \tilde{X}_{1:j-1})$ is the conditional distribution of X_j given $(X_{-j}, \tilde{X}_{1:j-1})$. When $n = 3$, this would work as follows: sample \tilde{X}_1 from $\mathcal{L}(X_1 | X_{2:3})$. Once this is done, $\mathcal{L}(X_{1:3}, \tilde{X}_1)$ is available and we, therefore, know $\mathcal{L}(X_2 | X_1, X_3, \tilde{X}_1)$. Hence, we can sample \tilde{X}_2 from this distribution. Continuing, $\mathcal{L}(X_{1:3}, \tilde{X}_{1:2})$ becomes known and we can sample \tilde{X}_3 from $\mathcal{L}(X_3 | X_{1:2}, \tilde{X}_{1:2})$.

It is not immediately clear why Algorithm 1 yields a sequence of random variables obeying the exchangeability property (3.1), and we prove this fact in Appendix A. There is, of course, nothing special about the ordering in which knockoffs are created and equally valid constructions may be obtained by looping through an arbitrary ordering of the variables. For example, in a data analysis application where we would need to build a knockoff copy for each row of the design, independent (random) orderings may be used.

To have power or, equivalently, to have a low Type II error rate, it is intuitive that we would like to have original features X_j and their knockoff companions \tilde{X}_j to be as “independent” as possible.

We do not mean to imply that running Algorithm 1 is a simple matter. In fact, it may prove rather complicated since we would have to recompute the conditional distribution at each step; this problem is left for future research. Instead, in this paper we shall work with approximate MF knockoffs and will demonstrate empirically that for models of interest, such constructions yield FDR control.

3.4.2 Approximate constructions: second-order model-free knockoffs

Rather than asking that $(X, \tilde{X})_{\text{swap}(S)}$ and (X, \tilde{X}) have the same distribution for any subset S , we can ask that they have the same first two moments, i.e., the same mean and covariance. Equality of means is a simple matter. As far as the covariances are concerned, equality is equivalent to

$$\text{cov}(X, \tilde{X}) = \mathbf{G}, \quad \text{where} \quad \mathbf{G} = \begin{bmatrix} \mathbf{\Sigma} & \mathbf{\Sigma} - \text{diag}\{s\} \\ \mathbf{\Sigma} - \text{diag}\{s\} & \mathbf{\Sigma} \end{bmatrix}. \quad (3.13)$$

We, of course, recognize the same form as in (3.2) where the parameter s is chosen to yield a positive semidefinite covariance matrix. (When (X, \tilde{X}) is Gaussian, a matching of the first two moments implies a matching of the joint distributions so that we have an exact construction.) Furthermore, Section 3.1.2 shows that the same problem was already solved in Barber and Candès (2015), as the same constraint on s applies but with the empirical covariance replacing the true covariance. This means that the same two constructions proposed in Barber and Candès (2015) are just as applicable to *second-order model-free knockoffs*.

For the remainder of this section, we will assume the covariates have each been translated and rescaled to have mean zero and variance one. To review, the *equicorrelated* construction uses

$$s_j^{\text{EQ}} = 2\lambda_{\min}(\mathbf{\Sigma}) \wedge 1 \text{ for all } j,$$

which minimizes the correlation between variable-knockoff pairs subject to the constraint that all such pairs must have the same correlation. The *semidefinite program (SDP)* construction solves the convex program

$$\begin{aligned} & \text{minimize} && \sum_j |1 - s_j^{\text{SDP}}| \\ & \text{subject to} && s_j^{\text{SDP}} \geq 0 \\ & && \text{diag}\{s^{\text{SDP}}\} \preceq 2\mathbf{\Sigma}, \end{aligned} \quad (3.14)$$

which minimizes the sum of absolute values of variable-knockoff correlations among all suitable s .

In applying these constructions to problems with large p , we run into some new difficulties:

- Excepting very specially-structured matrices like the identity, $\lambda_{\min}(\mathbf{\Sigma})$ tends to be extremely small as p gets large. The result is that constructing equicorrelated knockoffs in high dimensions, while fairly computationally easy, will result in very low power, since all the original variables will be nearly indistinguishable from their knockoff counterparts.
- For large p , (3.14), while convex, is prohibitively computationally expensive. However, if it could be computed, it would produce much larger s_j ’s than the equicorrelated construction and thus be considerably more powerful.

To address these difficulties, we first generalize the two knockoff constructions by the following two-step procedure, which we call the approximate semidefinite program (ASDP) construction:

Step 1. Choose an approximation $\mathbf{\Sigma}_{\text{approx}}$ of $\mathbf{\Sigma}$ and solve:

$$\begin{aligned} & \text{minimize} && \sum_j |1 - \hat{s}_j| \\ & \text{subject to} && \hat{s}_j \geq 0 \\ & && \text{diag}\{\hat{s}\} \preceq 2\mathbf{\Sigma}_{\text{approx}}. \end{aligned} \quad (3.15)$$

Algorithm 2 Conditional Randomization Test.

Input: A set of n independent samples $(X_{i1}, \dots, X_{ip}, Y_i)_{1 \leq i \leq n}$ assembled in a data matrix \mathbf{X} and a response vector y , a feature importance statistic $T_j(\mathbf{X}, y)$ to test whether X_j and Y are conditionally independent.

Loop: for $k = 1, 2, \dots, K$ do

Create a new data matrix $\mathbf{X}^{(k)}$ by simulating the j th column of \mathbf{X} from $\mathcal{L}(X_j | X_{-j})$ (and keeping the remaining columns the same). That is, $X_{ij}^{(k)}$ is sampled from the conditional distribution $X_{ij} | \{X_{i1}, \dots, X_{ip}\} \setminus \{X_{ij}\}$, and is (conditionally) independent of X_{ij} .

Output: A (one-sided) p-value

$$P_j = \frac{1}{K+1} \left[1 + \sum_{k=1}^K \mathbb{1}_{T_j(\mathbf{X}^{(k)}, y) \geq T_j(\mathbf{X}, y)} \right].$$

As with permutation tests, adding one in both the numerator and the denominator makes sure that the null p-values are stochastically larger than uniform variables.

Step 2. Solve:

$$\begin{aligned} & \text{maximize} && \gamma \\ & \text{subject to} && \text{diag}\{\gamma \hat{s}\} \preceq 2\mathbf{\Sigma}, \end{aligned} \tag{3.16}$$

and set $s^{\text{ASDP}} = \gamma \hat{s}$. Note this problem can be solved quickly by, e.g., bisection search over $\gamma \in [0, 1]$.

ASDP with $\mathbf{\Sigma}_{\text{approx}} = \mathbf{I}$ trivially gives $\hat{s}_j = 1$ and $\gamma = 2\lambda_{\min}(\mathbf{\Sigma}) \wedge 1$, reproducing the equicorrelated construction. ASDP with $\mathbf{\Sigma}_{\text{approx}} = \mathbf{\Sigma}$ clearly gives $\hat{s}_j = s^{\text{SDP}}$ and $\gamma = 1$, reproducing the SDP construction. Note that the ASDP step 2 is always fast, so the speed of the equicorrelated construction comes largely because the problem *separates* into p computationally independent SDP subproblems of $\min |1 - \hat{s}_j|$ s.t. $0 \leq \hat{s}_j \leq 2$. However, power is lost due to the very naïve approximation $\mathbf{\Sigma}_{\text{approx}} = \mathbf{I}$ which results in a very small γ .

In general, we can choose $\mathbf{\Sigma}_{\text{approx}}$ to be an m -block-diagonal approximation of $\mathbf{\Sigma}$, so that the ASDP from Step 1 separates into m smaller, more computationally tractable, and trivially parallelizable SDP subproblems. If the approximation is fairly accurate, we may also find that γ remains large, so that the knockoffs are nearly as powerful as if we had used the SDP construction. We demonstrate the ASDP construction in Section 6 when we analyze the Crohn’s disease data.

4 The conditional randomization test

This section presents an alternative approach to the controlled variable selection problem. To describe our approach, it may be best to consider an example. Assume we are in a regression setting and let $\hat{b}_j(\lambda)$ be the value of the Lasso estimate of the j th regression coefficient. We would like to use the statistic $\hat{b}_j(\lambda)$ to test whether Y is conditionally independent of X_j since large values of $|\hat{b}_j(\lambda)|$ provide evidence against the null. To construct a test, however, we would need to know the sampling distribution of $\hat{b}_j(\lambda)$ under the null hypothesis that Y and X_j are conditionally independent, and it is quite unclear how one would obtain such knowledge.

4.1 The test

A way out is to sample the covariate X_j conditional on all the other covariates (but not the response), where by “sample” we explicitly mean to draw a new sample from the conditional distribution of $X_j | X_{-j}$ using a random number generator. We then compute the Lasso statistic $\hat{b}_j^*(\lambda)$, where the $*$ superscript indicates that the statistic is computed from the artificially sampled value of the covariate X_j . Now, under the null hypothesis of conditional independence between Y and X_j , it happens that $\hat{b}_j^*(\lambda)$ and $\hat{b}_j(\lambda)$ are identically distributed and that, furthermore, this statement holds true conditional on Y and all the other covariates. This claim is proved in Lemma 4.1 below. A consequence of this is that by simulating a covariate conditional on the others, we can sample at will from the conditional distribution of any test statistic and compute p-values as described in Algorithm 2.

Lemma 4.1. *Let (Z_1, Z_2, Y) be a triple of random variables, and construct another triple (Z_1^*, Z_2, Y) as*

$$Z_1^* | (Z_2, Y) \stackrel{d}{=} Z_1 | Z_2.$$

Then under the null hypothesis $Y \perp\!\!\!\perp Z_1 \mid Z_2$, any test statistic $T = t(Z_1, Z_2, Y)$ obeys

$$T \mid (Z_2, Y) \stackrel{d}{=} T^* \mid (Z_2, Y),$$

where $T^* = t(Z_1^*, Z_2, Y)$.

Proof. To prove the claim, it suffices to show that Z_1 and Z_1^* have the same distribution conditionally on (Z_2, Y) . This follows from

$$Z_1^* \mid (Y, Z_2) \stackrel{d}{=} Z_1 \mid Z_2 \stackrel{d}{=} Z_1 \mid (Z_2, Y).$$

The first equality comes from the definition of Z_1^* while the second follows from the conditional independence of Y and Z_1 , which holds under the null. \square

The consequence of Lemma 4.1 is that we can compute the 95% percentile, say, of the conditional distribution of T^* denoted by $t_{0.95}^*(Z_2, Y)$. Then by definition, under the null,

$$\mathbb{P}(T > t_{0.95}^*(Z_2, Y) \mid (Z_2, Y)) \leq 0.05.$$

Since this equality holds conditionally, it also holds marginally.

4.2 Literature review

The conditional randomization test is most closely related to the propensity score (Rosenbaum and Rubin, 1983), which also uses the conditional distribution $X_j \mid X_{-j}$ to perform inference on the conditional relationship between Y and X_j given X_{-j} . However, propensity scores require X_j be binary, and the propensity score itself is normally estimated, although Rosenbaum (1984) shows that when all the covariates jointly take a small number of discrete values, propensity score analysis can be done exactly. Doran et al. (2014) also rely on the data containing repeated observations of X_{-j} so that certain observations can be permuted nonparametrically while maintaining the null distribution. In fact, the exact term “conditional randomization test” has also been used in randomized controlled experiments to test for independence of Y and X_j conditioned more generally on some function of X_{-j} (such as a measure of imbalance in X_{-j} if X_j is binary), again relying on discreteness of the function so that there exist permutations of X_j which leave the function value unchanged. Despite the similar name, our conditional randomization test is quite distinct from these, as it does not rely on discreteness or experimental control in any of the covariates.

Another line of work exists within the linear model regime, whereby the null (without X_j) model is estimated and then the empirical residuals are permuted to produce a null distribution for Y (Freedman and Lane, 1983). Because this approach is only exact when the empirical residuals match the true residuals, it explicitly relies on a parametric model for $Y \mid X_{-j}$, as well as the ability to estimate it quite accurately.

4.3 Comparisons with knockoffs

One major limitation of the conditional randomization method is its computational cost. It requires computing randomization p-values for many covariates and to a high-enough resolution for multiple-comparisons correction. Clearly, this requires samples in the extreme tail of the p-value distribution. This means computing a very large number of feature importance statistics T_j , each of which can be expensive since for reasons outlined in the drawbacks associated with marginal testing, powerful T_j ’s will take into account the full dimensionality of the model, e.g., absolute value of the Lasso-estimated coefficient. In fact, the number of computations of T_j , tallied over all j , required by the conditional randomization method is $\Omega(p)$.⁷ To see this, suppose for simplicity that all R rejected p-values take on the value of half the BHq cutoff equal to $\tau = qR/p$, and all we need to do is upper-bound them below τ . This means there are R p-values P_j for which plugging $K = \infty$ into Algorithm 2 would yield $P_j = \tau/2$. After $K < \infty$ samples, the approximate p-value (ignoring the +1 correction) is distributed as $K^{-1} \text{Bin}(K, P_j)$. We could then use this binomial count to construct a confidence interval for P_j . A simple calculation shows that to be reasonably confident that $P_j \leq \tau$, K must be on the order of at least $1/\tau$. Since there are R such p-values, this justifies the claim.

Note that for knockoffs, the analogous computation of T need only be done exactly once. If, for instance, each T_j requires a Lasso computation, then the conditional randomization test’s computational burden is very challenging for medium-scale p in the thousands and prohibitive for large-scale (e.g., genetics) p in the hundreds of thousands or millions. We will see in Section 5.3.1 that there are power gains, along with huge computational costs, to be had

⁷ $a(N) \in \Omega(b(N))$ means that there exist N_0 and $C > 0$ such that $a(N) \geq Cb(N)$ for $N \geq N_0$.

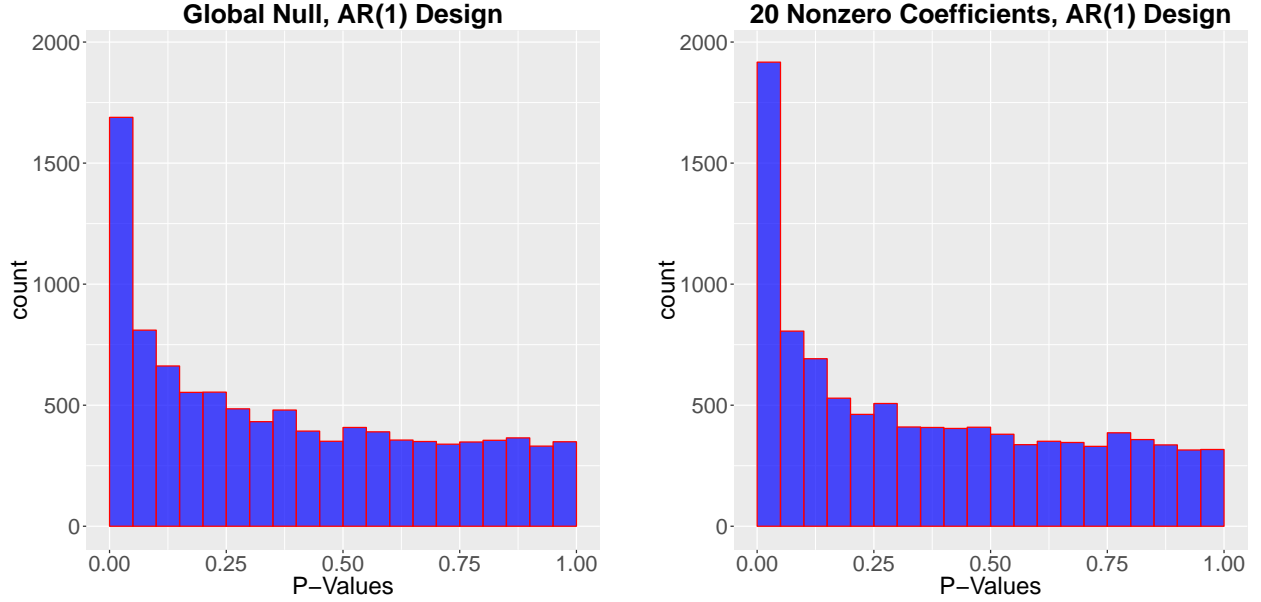


Figure 1: Distribution of null logistic regression p-values with $n = 500$ and $p = 200$; 10,000 replications.

by using conditional randomization in place of knockoffs, and Section 6 will show that the MF knockoff procedure easily scales to large data sets.

Another advantage of MF knockoffs is its guaranteed control of the FDR, whereas the BHq procedure does not offer strict control when applied to arbitrarily dependent p-values.

5 Numerical simulations

In this section we demonstrate the importance, utility, and practicality of MF knockoffs for high-dimensional non-parametric conditional modeling.

5.1 Logistic regression p-values

Asymptotic maximum likelihood theory promises valid p-values for each coefficient in a GLM only when $n \gg p$. However, these approximate p-values can usually be computed as long as $n > p$, so a natural question arising from high-dimensional applications is whether such asymptotic p-values are valid when n and p are both large with $p/n \geq 0.1$, for example. We simulated 10^4 independent design matrices ($n = 500$, $p = 200$) and binary responses from a logistic regression for the following two settings:

- (1) (X_1, \dots, X_p) is an AR(1) time series with AR coefficient 0.5 and

$$Y \mid X_1, \dots, X_p \sim \text{Bernoulli}(0.5)$$

- (2) (X_1, \dots, X_p) is an AR(1) time series with AR coefficient 0.5 and

$$Y \mid X_1, \dots, X_p \sim \text{Bernoulli}(\text{logit}(0.08(X_2 + \dots + X_{21})))$$

Histograms for the p-values for β_1 (null in all cases) are shown in Figure 1. Both histograms are far from uniform, and Table 1 shows each distribution's concentration near zero. We see that the small quantiles have extremely inflated probabilities—over 20 times nominal for $\mathbb{P}\{p\text{-value} \leq 0.1\%\}$ in setting (2). We also see that the exact null distribution depends on the unknown coefficient sequence β_2, \dots, β_p , since the probabilities between settings differ statistically significantly at all three cutoffs.

To confirm that this non-uniformity is not just a finite-sample effect, we also simulated 10^4 i.i.d. $\mathcal{N}(0, 1)$ design matrices with independent Bernoulli(0.5) responses for $n = 500$, $p = 200$ and $n = 5000$, $p = 2000$ as settings (3) and (4), respectively. Table 1 shows that the distribution does not really change as n and p are increased with constant proportion.

	(1)	(2)	(3)	(4)
$\mathbb{P}\{p\text{-value} \leq 5\%\}$	16.89% (0.37%)	19.17% (0.39%)	16.88% (0.37%)	16.78% (0.37%)
$\mathbb{P}\{p\text{-value} \leq 1\%\}$	6.78% (0.25%)	8.49% (0.28%)	7.02% (0.26%)	7.03% (0.26%)
$\mathbb{P}\{p\text{-value} \leq 0.1\%\}$	1.53% (0.12%)	2.27% (0.15%)	1.87% (0.14%)	2.04% (0.14%)

Table 1: Inflated p-value probabilities with estimated Monte Carlo standard errors in parentheses. See text for meanings of settings (1), (2), (3), (4).

These results show that the usual logistic regression p-values one might use when $n \geq p$ can have null distributions that are quite far from uniform, and even if one wanted to correct that distribution, it depends in general on unknown problem parameters, further complicating matters. When $n < p$ the problem becomes even more challenging, with existing methods similarly asymptotic as well as requiring stringent sparsity assumptions (van de Geer et al., 2014). Thus, despite the wealth of research on controlling FDR, without a way to obtain valid p-values, even the problem of controlling FDR in medium-to-high-dimensional GLMs remains unsolved.

5.2 Alternative knockoff statistics

As mentioned in Section 3.2, the new MF knockoff framework allows for a wider variety of W statistics to be used than in the original knockoffs framework. Choices of Z_j include well-studied statistical measures such as the coefficient estimated in a GLM, but can also include much more ad-hoc/heuristic measures such as random forest bagging feature importances or sensitivity analysis measures such as the Monte-Carlo-estimated total sensitivity index. By providing model selection with rigorous Type I error control for general models and statistics, knockoffs can be used to improve the interpretability of complex black-box supervised/machine learning models. There are also many available choices for the anti-symmetric function f_j , such as $|Z_j| - |\tilde{Z}_j|$, $\text{sign}(|Z_j| - |\tilde{Z}_j|) \max\{|Z_j|, |\tilde{Z}_j|\}$, or $\log(|Z_j|) - \log(|\tilde{Z}_j|)$.

We discuss here a few appealing new options for statistics W , but defer full exploration of these very extensive possibilities to future work.

5.2.1 Adaptive knockoff statistics

The default statistic suggested in Barber and Candès (2015) is the Lasso Signed Max (LSM), which corresponds to Z_j being the largest penalty parameter at which the j th variable enters the model in the Lasso regression of y on $[\mathbf{X}, \tilde{\mathbf{X}}]$, and $f_j = \text{sign}(|Z_j| - |\tilde{Z}_j|) \max\{|Z_j|, |\tilde{Z}_j|\}$. In addition to the LSM statistic, Barber and Candès (2015) suggested alternatives such as the difference in absolute values of estimated coefficients for a variable and its knockoff:

$$W_j = |\hat{b}_j| - |\hat{b}_{j+p}|,$$

where the \hat{b}_j, \hat{b}_{j+p} are estimated so that W obeys the sufficiency property required by the original knockoff procedure, e.g., by ordinary least squares or the Lasso with a pre-specified tuning parameter. The removal of the sufficiency requirement for MF knockoffs allows us to improve this class of statistics by adaptively tuning the fitted model. The simplest example is the LCD statistic introduced in Section 3.2, which uses cross-validation to choose the tuning parameter in the Lasso. Note the LCD statistic can be easily extended to any GLM by replacing the first term in (3.7) by a non-Gaussian negative log-likelihood, such as in logistic regression; we will refer to all such statistics generically as LCD. The key is that the tuning and cross-validation is done on the augmented design matrix $[\mathbf{X}, \tilde{\mathbf{X}}]$, so that W still obeys the flip-sign property.

More generally, MF knockoffs allows us to construct statistics that are highly adaptive to the data, as long as that adaptivity does not distinguish between original and knockoff variables. For instance, we could compute the cross-validated error of the ordinary Lasso (still of y on $[\mathbf{X}, \tilde{\mathbf{X}}]$) and compare it to that of a random forest, and choose Z to be a feature importance measure derived from whichever one has smaller error. Since the Lasso works best when the true model is close to linear, while random forests work best in non-smooth models, this approach gives us high-level adaptivity to the model smoothness, while the MF knockoff framework ensures strict Type I error control.

Returning to the simpler example of adaptivity, we found the LCD statistic to be uniformly more powerful than the LSM statistic across a wide range of simulations (linear and binomial GLMs, ranging covariate dependence, effect size, sparsity, sample size, total number of variables), particularly under covariate dependence. We note, however, the importance of choosing the penalty parameter that minimizes the cross-validated error, as opposed to the default

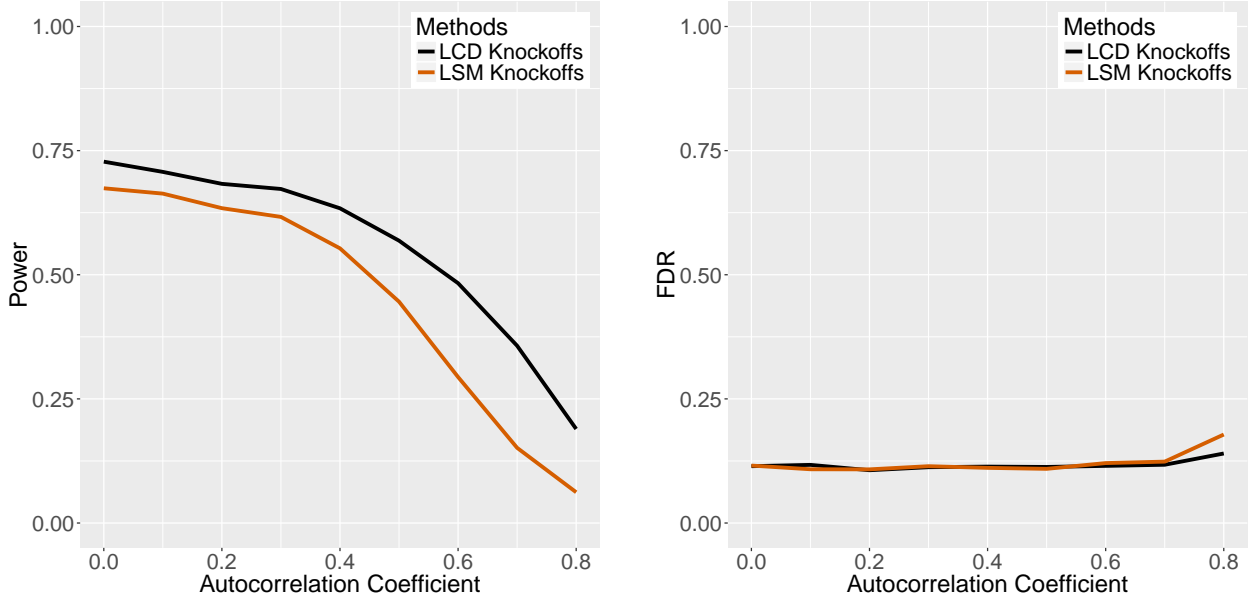


Figure 2: Power and FDR (target is 10%) for knockoffs with the LCD and LSM statistics. The design matrix has i.i.d. rows and AR(1) columns with autocorrelation coefficient specified by the x-axes of the plots, and marginally each $X_j \sim \mathcal{N}(0, 1/n)$. Here, $n = 3000$, $p = 1000$, and y comes from a binomial linear model with logit link function with 60 nonzero regression coefficients of magnitude 3.5 and random signs. Each point represents 200 replications.

in some computational packages of using the “one standard error” rule, as the latter causes LCD to be underpowered compared to LSM in low-power settings. Figure 2 shows a simulation with $n = 3000$, $p = 1000$ of a binomial linear model (with statistics computed from Lasso logistic regression) that is representative of the power difference between the two statistics. In all our simulations, unless otherwise specified, MF knockoffs is always run using the LCD statistic. Explicitly, when the response variable is continuous, we use the standard Lasso with Gaussian linear model likelihood, and when the response is binary, we use Lasso-penalized logistic regression.

5.2.2 Bayesian knockoff statistics

Another very interesting source of knockoff statistics comes from Bayesian procedures. If a statistician has prior knowledge about the problem, he or she can encode it in a Bayesian model and use the resulting estimators to construct a statistic (e.g., difference of absolute posterior mean coefficients, or difference or log ratio of posterior probabilities of nonzero coefficients with a sparse prior). What makes this especially appealing is that the statistician gets the power advantages of incorporating prior information, while maintaining a strict frequentist guarantee on the Type I error, *even if the prior is false!*

As an example, we ran knockoffs in an experiment with a Bayesian hierarchical regression model with $n = 300$, $p = 1000$, and $\mathbb{E}(\|\beta\|_0) = 60$ ($\|\cdot\|_0$ denotes the ℓ_0 norm, or the number of nonzero entries in a vector); see Appendix B for details. We chose a simple canonical model with Gaussian response to demonstrate our point, but the same principle applies to more complex, nonlinear, and non-Gaussian Bayesian models as well. The statistics we used were the LCD and a Bayesian variable selection (BVS) statistic, namely $Z_j - \tilde{Z}_j$ where Z_j and \tilde{Z}_j are the posterior probabilities that the j th original and knockoff coefficients are nonzero, respectively (George and McCulloch, 1997); again see Appendix B for details. Figure 3 shows that the accurate prior information supplied to the Bayesian knockoff statistic gives it improved power over LCD which lacks such information, but that they have the same FDR control (and they would even if the prior information were incorrect).

5.3 Alternative procedures

To assess the relative power of knockoffs, we compare to a number of alternatives in settings in which they are valid:

1. The original knockoff procedure with settings recommended in Barber and Candès (2015). This method can only be applied in homoscedastic Gaussian linear regression when $n \geq p$.

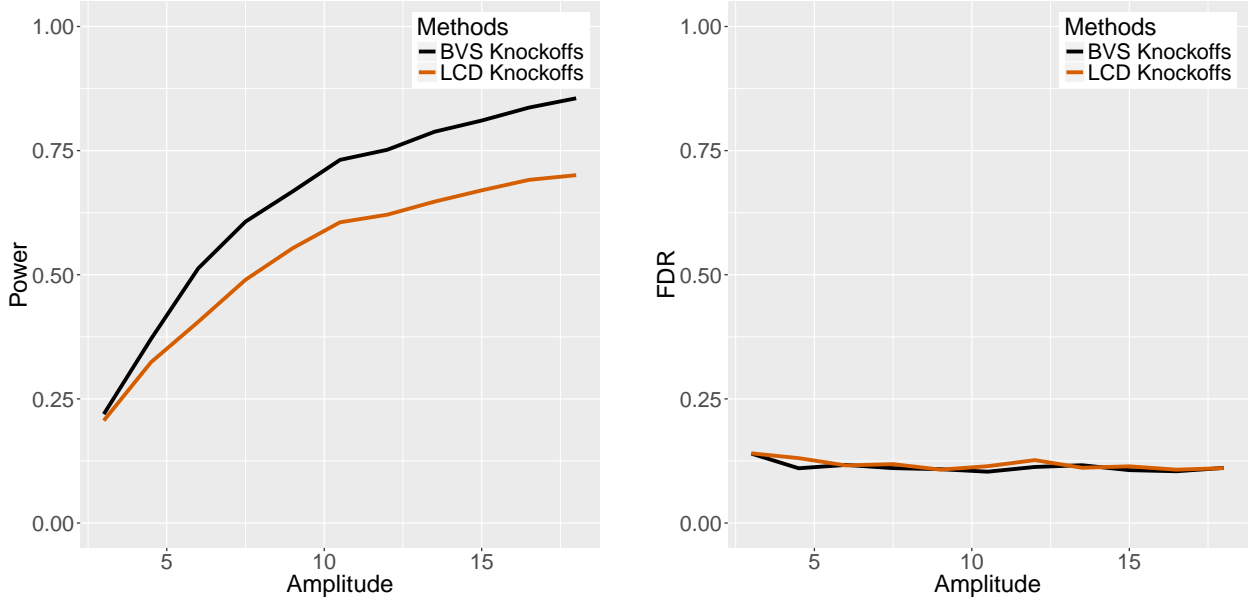


Figure 3: Power and FDR (target is 10%) for knockoffs with the LCD and BVS statistics. The design matrix is i.i.d. $\mathcal{N}(0, 1/n)$, $n = 300$, $p = 1000$, and y comes from a Gaussian linear model with β and the noise variance randomly chosen (see Appendix B for the precise model). Here, the nonzero entries of β are Gaussian with mean zero and standard deviation given on the x-axis; the expected number of nonzero components is 60; the expected variance of the noise is 1. Each point represents 200 replications.

2. BHq applied to asymptotic GLM p-values. This method can only be applied when $n \geq p$, and although for linear regression exact p-values can be computed (when the MLE exists), for any other GLM these p-values can be far from valid unless $n \gg p$, as shown in Section 5.1.
3. BHq applied to marginal test p-values. The correlation between the response and each covariate is computed and compared to its null distribution, which under certain Gaussian assumptions is closed-form, but in general can at least be simulated exactly by conditioning on y and using the known marginal distribution of X_j . Although these tests are valid for testing hypotheses of *marginal* independence (regardless of n and p), such hypotheses only agree with the desired *conditional* independence hypotheses when the covariates are exactly independent of one another.
4. BHq applied to the p-values from the conditional randomization test described in Section 4.

Recall from the drawbacks to marginal testing in Section 1.5 that, even ignoring the marginal validity of the p-values in procedures 2–4, the joint distribution of the p-values will not in general satisfy the assumptions for BHq to control FDR. However, the conservative alternative had extremely noncompetitive power in every simulation we tried, so we only report results for BHq.

5.3.1 Comparison with conditional randomization

We start by comparing MF knockoffs with procedure 4, BHq applied to conditional randomization test p-values, for computational reasons. We simulated $n = 400$ i.i.d. rows of $p = 600$ AR(1) covariates with autocorrelation 0.3, and response following a logistic regression model with 40 nonzero coefficients of random signs. Figure 4 shows the power and FDR curves as the coefficient amplitude was varied. We see that the conditional randomization test gives higher power with similar FDR control, but this comes at a hugely increased computational cost. This simulation has considerably smaller n and p than any other simulation in the paper, and we still had to apply a number of computational speed-ups/shortcuts, described below, to keep the computation time within reason.

- As mentioned, the LCD statistic is a powerful one for MF knockoffs, so we wanted to compare it to the analogue for the conditional randomization test, namely, the absolute value of the Lasso-estimated coefficient. However, choosing the penalty parameter by cross-validation turned out to be too expensive, so instead we chose a fixed value by simulating repeatedly from the known model (with amplitude 30), running cross-validation for each

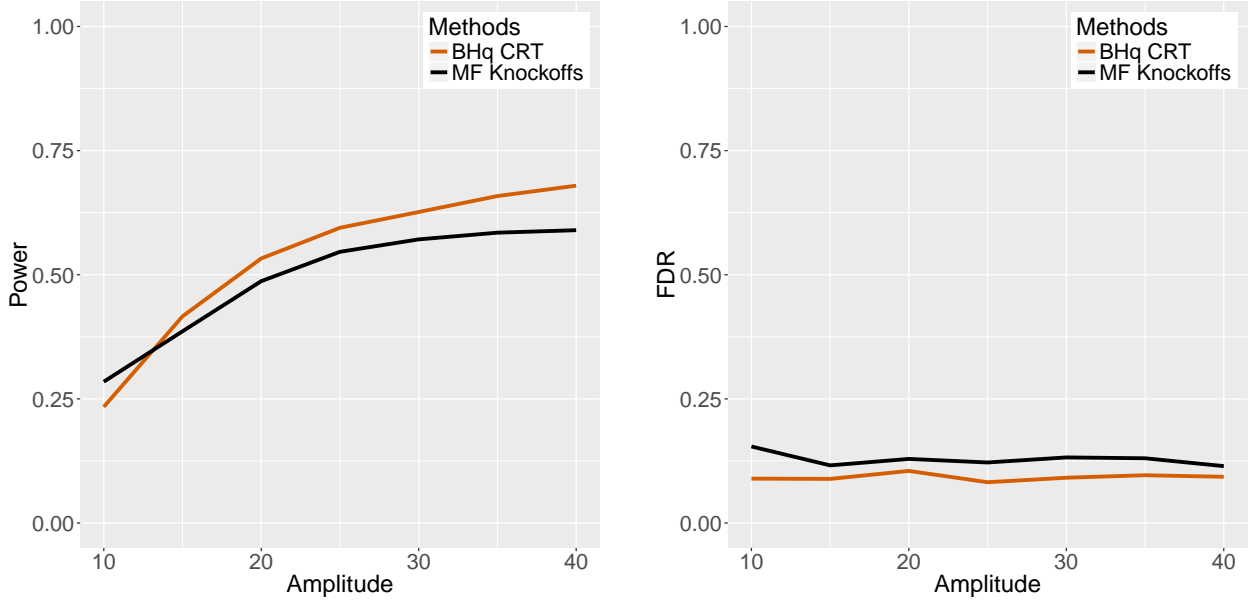


Figure 4: Power and FDR (target is 10%) for MF knockoffs and BHq applied to conditional randomization test p-values. The design matrix has i.i.d. rows and AR(1) columns with autocorrelation 0.3, $n = 400$, $p = 600$, and y comes from a binomial linear model with logit link function with $\|\beta\|_0 = 40$, and all nonzero entries of β having equal magnitudes and random signs; each point represents 200 replications.

repetition, and choosing the median of the cross-validation-error-minimizing penalty parameters (the chosen value was 0.00053). For a fair comparison, we did the same for MF knockoffs (except the simulations included knockoff variables too—the chosen value was 0.00077). This speed-up is of course impossible in practice, as the true model is not known. It is not clear how one would choose the penalty parameter if not by cross-validation (at considerable extra computational expense), although a topic of current research is how to choose tuning parameters efficiently and without explicit reliance on a model for the response, e.g., (Lv and Liu, 2014).

- Because the statistic used was the (absolute value of the) estimated coefficient from a sparse estimation procedure, many of the observed statistics were exactly zero, and for these, the conditional randomization p-values can be set to one without further computation. This did not require any prior knowledge, although it will only work for statistics whose distribution has a point mass at the smallest possible value.
- Because the power was calibrated, we knew to expect at least around 10 discoveries, and thus could anticipate the BHq cutoff being at least $0.1 \times 10/600$. This cutoff gives a sense of the p-value resolution needed, and we chose the number of randomizations to be roughly 10 times the inverse of the BHq cutoff, namely, 10,000. However, we made sure that all 10,000 randomizations were only computed for very few covariates, both using the previous bullet point and also by checking after periodic numbers of randomizations whether we can reject the hypothesis that the true p-value is below the approximate BHq cutoff upper-bound of $0.1 \times 44/600$ (the 44 comes from 40 nonzeros with 10% false discoveries). For instance, after just 10 randomizations, if the approximate p-value so far is greater than 0.2, we can reject the hypothesis that the exact p-value is below $0.1 \times 44/600$ at significance 0.0001 (and thus compute no further randomizations for that covariate). Speed-ups like this are possible in practice, although they require having an idea of how many rejections will be made, which is not generally available ahead of time.

With these speed-ups, Figure 4 took roughly three years of serial computation time, while the MF knockoffs component took only about six hours, or about 5000 times less (all computation was run in Matlab 2015b, and both methods used glmnet to compute statistics). Because of the heavy computational burden, we were unable to include the conditional randomization test in our further, larger simulations—recall from Section 4.3 that the number of T_j computations scales optimistically linearly in p . To summarize, conditional randomization testing appears somewhat more powerful than MF knockoffs, but is computationally infeasible for large data sets (like that in Section 6).

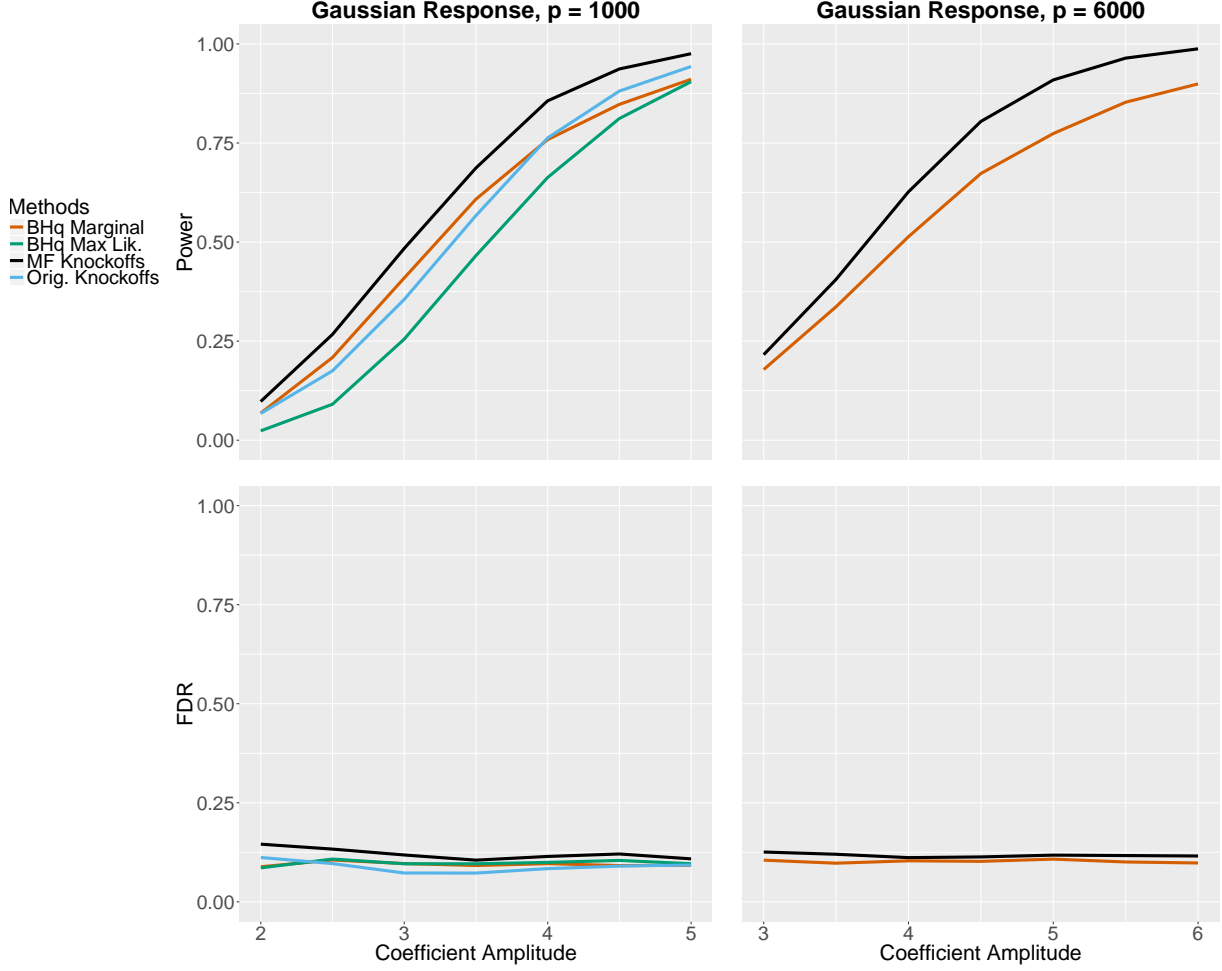


Figure 5: Power and FDR (target is 10%) for MF knockoffs and alternative procedures. The design matrix is i.i.d. $\mathcal{N}(0, 1/n)$, $n = 3000$, $p =$ (left): 1000 and (right): 6000, and y comes from a Gaussian linear model with 60 nonzero regression coefficients having equal magnitudes and random signs. The noise variance is 1. Each point represents 200 replications.

5.3.2 Effect of signal amplitude

Our first simulation comparing MF knockoffs to procedures 1–3 is by necessity in a Gaussian linear model with $n > p$ and independent covariates—the only setting in which all procedures approximately control the FDR. Specifically, the left side of Figure 5 plots the power and FDR for the four procedures when $X_{ij} \stackrel{i.i.d.}{\sim} \mathcal{N}(0, 1/n)$, $n = 3000$, $p = 1000$, $\|\beta\|_0 = 60$, the noise variance $\sigma^2 = 1$, and the nonzero entries of β have random signs and equal magnitudes, varied along the x-axis. All methods indeed control the FDR, and MF knockoffs is the most powerful, with as much as 10% higher power than its *nearest* alternative. The right side of Figure 5 shows the same setup but in high dimensions: $p = 6000$. In the high-dimensional regime, neither maximum likelihood p-values nor original knockoffs can even be computed, and the MF knockoff procedure has considerably higher power than BHq applied to marginal p-values.

Next we move beyond the Gaussian linear model to a binomial linear model with logit link function, precluding the use of the original knockoff procedure. Figure 6 shows the same simulations as Figure 5 but with Y following the binomial model. The results are similar to those for the Gaussian linear model, except that BHq applied to the asymptotic maximum likelihood p-values now has an FDR above 50% (rendering its high power meaningless), which can be understood as a manifestation of the phenomenon from Section 5.1. In summary, MF knockoffs continues to have the highest power among FDR-controlling procedures.

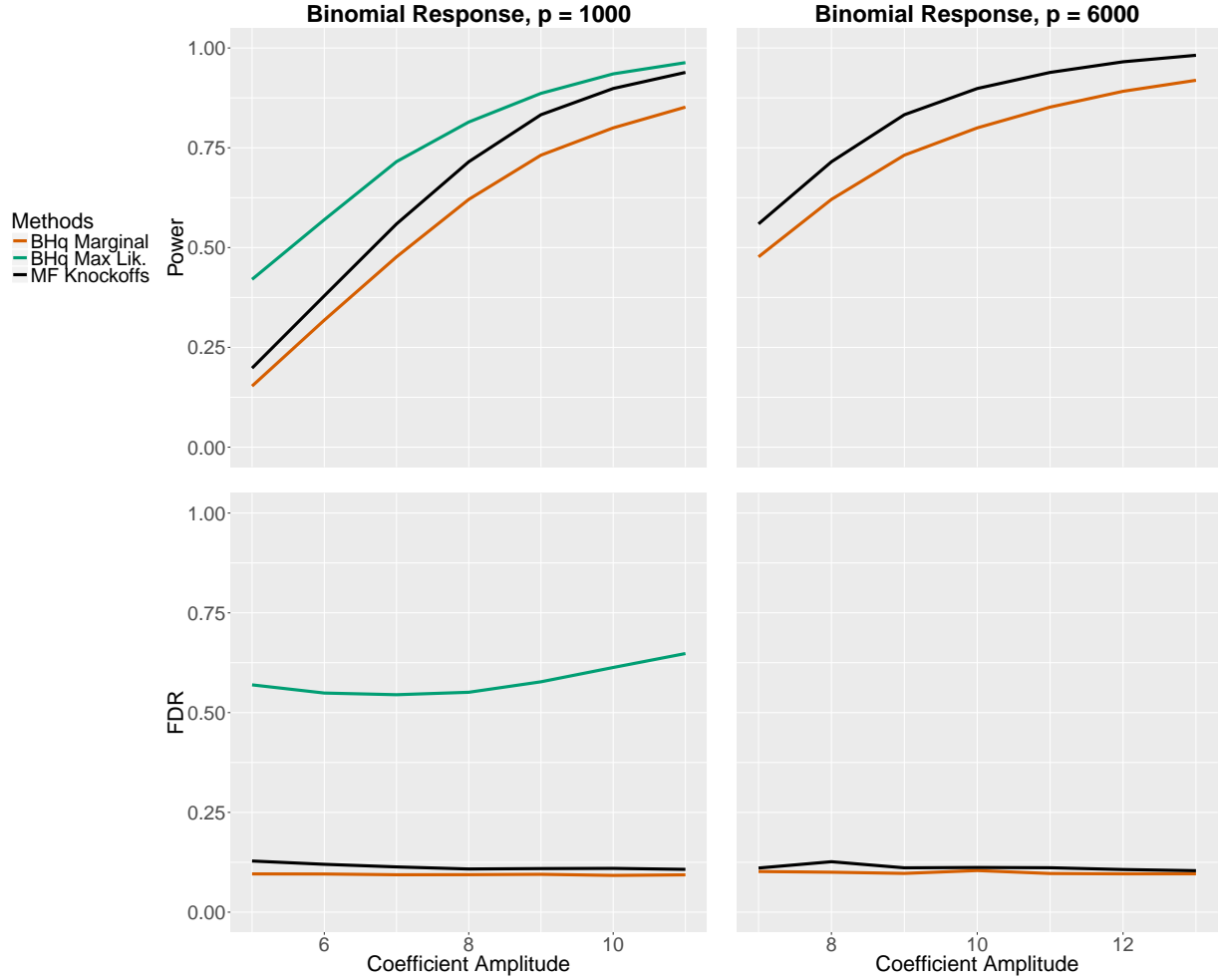


Figure 6: Power and FDR (target is 10%) for MF knockoffs and alternative procedures. The design matrix is i.i.d. $\mathcal{N}(0, 1/n)$, $n = 3000$, $p =$ (left): 1000 and (right): 6000, and y comes from a binomial linear model with logit link function, and 60 nonzero regression coefficients having equal magnitudes and random signs. Each point represents 200 replications.

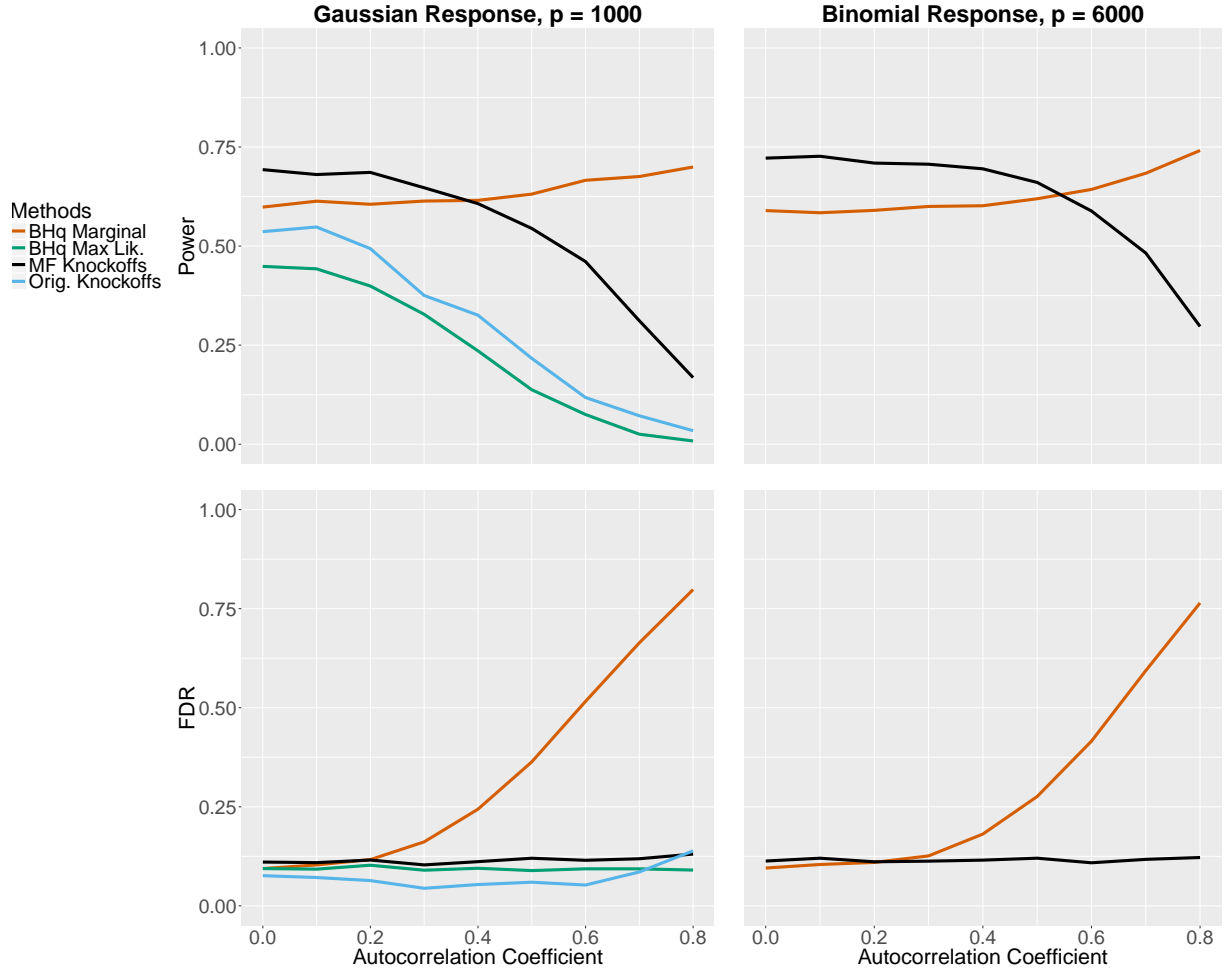


Figure 7: Power and FDR (target is 10%) for MF knockoffs and alternative procedures. The design matrix has i.i.d. rows and AR(1) columns with autocorrelation coefficient specified by the x-axes of the plots, and marginally each $X_j \sim \mathcal{N}(0, 1/n)$. (left): $n = 3000$, $p = 1000$, and y follows a Gaussian linear model. (right): $n = 3000$, $p = 6000$, and y follows a binomial linear model with logit link function. In both cases, there are 60 nonzero coefficients having magnitudes equal to 3.5 on the left and 10 on the right, random signs, and randomly selected locations. Each point represents 200 replications.

5.3.3 Effect of covariate dependence

To assess the relative power and FDR control of MF knockoffs as a function of covariate dependence, we ran similar simulations as in the previous section, but with covariates that are AR(1) with varying autocorrelation coefficient (while the coefficient amplitude remains fixed). It is now relevant to specify that the locations of the nonzero coefficients are uniformly distributed on $\{1, \dots, p\}$. In the interest of space, we only show the low-dimensional ($p = 1000$) Gaussian setting (where all four procedures can be computed) and the high-dimensional ($p = 6000$) binomial setting, as little new information is contained in the plots for the remaining two settings. Figure 7 shows that, as expected, BHq with marginal testing quickly loses FDR control with increasing covariate dependence. This is because the marginal tests are testing the null hypothesis of *marginal* independence between covariate and response, while recall from Definition 2.1 that all conditionally independent covariates are considered null, even if they are marginally dependent on the response. Concentrating on the remaining methods and just the left-hand part of the BHq marginal curves where FDR is controlled, Figure 7 shows that MF knockoffs continues to be considerably more powerful than alternatives as covariate dependence is introduced, in low- and high-dimensional linear and nonlinear models.

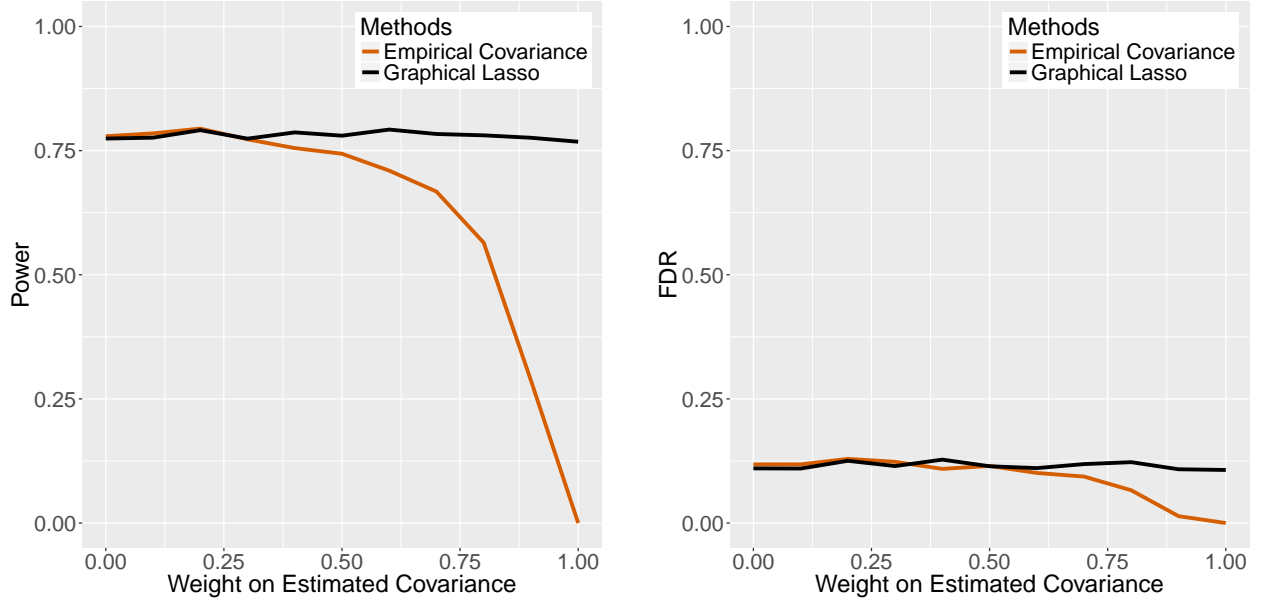


Figure 8: Power and FDR (target is 10%) for knockoffs with the LCD statistic as the covariance matrix used to generate knockoffs ranges from the truth to an estimated covariance; see text for details. The design matrix has i.i.d. rows and AR(1) columns with autocorrelation coefficient 0.3, and the matrix (including knockoffs) is standardized so that each column has mean zero and Euclidean norm 1. Here, $n = 800$, $p = 1500$, and y comes from a binomial linear model with logit link function with 50 nonzero entries having magnitude 20 and random signs. Each point represents 200 replications.

5.4 Robustness to overfitting

In many real applications, the true joint covariate distribution may not be known exactly, forcing the user to estimate it from the available data. As already mentioned, this is a challenging problem by itself, but often we have considerable outside information or unsupervised data that can be brought to bear to improve estimation. This raises the important question of how robust MF knockoffs is to error in the joint covariate distribution, particularly error that biases that distribution toward the empirical covariate distribution, often referred to as overfitting error. To answer this question, we generated knockoffs for Gaussian variables, but instead of using the true covariance matrix, we used a combination of the true covariance matrix and an estimated covariance matrix. Figure 8 shows the power and FDR as the covariance we use ranges from the true covariance matrix (AR(1) with autocorrelation 0.3) to an estimate based on the data only. The curve labeled Empirical Covariance represents knockoffs run using the following convex combination of the true and empirical covariance matrices:

$$\Sigma_{\text{EC}} = (1 - \alpha)\Sigma + \alpha\hat{\Sigma},$$

where α is the value plotted on the x-axis (Weight on Estimated Covariance) and $\hat{\Sigma}$ is the empirical covariance matrix. The curve labeled Graphical Lasso represents knockoffs run using the following combination of the true and empirical precision matrices:

$$\Sigma^{\text{temp}} = \left((1 - \alpha)\Sigma^{-1} + \alpha\hat{\Theta} \right)^{-1}, \quad \Sigma_{\text{GL}} = \text{diag}(r) \Sigma^{\text{temp}} \text{diag}(r)$$

where α is again the value plotted on the x-axis, the rescaling vector r has $r_j = \sqrt{\Sigma_{jj}/\Sigma_{jj}^{\text{temp}}}$, and $\hat{\Theta}$ is the inverse covariance estimated by the graphical Lasso (with penalty parameter chosen by 2-fold cross-validation on the log-likelihood).

The graphical Lasso is well suited for this problem since the covariates have a sparse precision matrix. When used, we see no change in either power or FDR as α varies from zero to one. What is more remarkable, then, is that the curve for the empirical covariance—a poor estimate of the true covariance given the high dimensionality—also maintains FDR control throughout the entire range! In fact, even when half of the covariance used is the empirical covariance, both power and FDR remain essentially the same as if the exact covariance matrix were used. Beyond

this point, MF knockoffs actually becomes conservative, with power and FDR approaching zero as the estimated covariance approaches the empirical covariance. This behavior at $\alpha = 1$ is not surprising, since $p > n$ and thus the empirical covariance is rank-deficient, forcing the knockoff variables to be exact replicas of their original counterparts.⁸ Intuitively, in the middle of the range, instead of treating the variables as coming from their true joint distribution, MF knockoffs treats them as coming from their true distribution “conditional” on being similar to their observed values.

6 Genetic analysis of Crohn’s disease

To test the robustness and practicality of the new knockoff procedure, we applied it to a data set containing genetic information on cases and controls for Crohn’s disease (CD). The data is provided by the Wellcome Trust Case Control Consortium and has been studied previously (WTCCC, 2007). After preprocessing (see Appendix C), there were $p = 377,749$ single nucleotide polymorphisms (SNPs) measured on $n = 4,913$ subjects (1,917 CD patients and 2,996 healthy controls). Although $p \gg n$, the joint dependence of SNPs has a strong spatial structure, and outside data can be used to improve estimation. In particular, we approximated the standardized joint distribution as multivariate Gaussian with covariance matrix estimated using the methodology of Wen and Stephens (2010), which shrinks the off-diagonal entries of the empirical covariance matrix using genetic distance information estimated from the HapMap CEU population. This approximation was used on each chromosome, and SNPs on different chromosomes were assumed to be independent. The statistic we use is the LCD. Although the data itself cannot be made available, all code is available at http://statweb.stanford.edu/~candes/MF_Knockoffs/.

One aspect of SNP data is that it contains some very high correlations, which presents two challenges to our methodology. The first is generic to the variable selection problem: it is very hard to choose between two or more nearly-identical (highly-correlated) variables if the data supports at least one of them being selected.⁹ To alleviate this, we clustered the SNPs using the estimated correlations as a similarity measure with a single-linkage cutoff of 0.5, and settle for discovering important SNP clusters. To do so we choose one representative from each cluster and approximate the null hypothesis that a cluster is conditionally independent of the response given the other clusters by the null hypothesis that a cluster *representative* is conditionally independent of the response given the other cluster *representatives*. To choose the representatives, we could ignore the observed data altogether and do something like pick representatives with the highest minor allele frequency (computed from outside data), and then run knockoffs as described in the paper. Although this produces a powerful procedure (about 50% more powerful than the original analysis by WTCCC (2007)), a more powerful approach is to select cluster representatives using a fraction of the observations, including their responses, such as by marginal testing. Note that such a data-splitting approach appears to make our null hypotheses random, as in the work on inference after selection reviewed in Section 1.4. However, the approximation we are making is that each representative stands for its cluster, and each cluster has exactly one associated null hypothesis, no matter how selection is performed, even if it were nonrandom. That is, the *approximate* hypotheses being tested do not actually depend on the selection (unlike Brzyski et al. (2016) where clusters are selected and where the very definition of a cluster actually depends on the selection), and our approach remains model-free, which together should make it clear that it is still quite different from the literature on inference after selection.

Explicitly, we randomly chose 20% of our observations and on those observations only, we ran marginal t-tests between each SNP and the response. Then from each cluster we chose the SNP with smallest t-test p-value to be the single representative of that cluster. Because the observations used for selecting cluster representatives have had their representative covariate values selected for dependence on the outcome, if we constructed knockoff variables as usual and included them in our procedure, the required exchangeability established in Lemma 3.2 would be violated. However, taking a page from Barber and Candès (2016), we can still use these observations in our procedure by making their knockoff variables just identical copies of the original variables (just for these observations). It is easy to show (see Appendix D) that constructing knockoffs in this way, as exact copies for the observations used to pick representatives and as usual for the remaining observations, the pairwise exchangeability of null covariates with their knockoffs is maintained. Of course, the observations with identical original and knockoff covariate values do not directly help the procedure distinguish between the original and knockoff variables, but including them improves the accuracy of the fitting procedure used to produce the feature importance statistics, so that power is improved indirectly because the Z_j become more accurate measures of feature importance.

Replacing clusters with single representatives reduces p down to 71,145 (so the average cluster size was just over

⁸Although in principle we could break ties and assign signs by coin flips when $W_j = 0$, we prefer to only select X_j with $W_j > 0$, as $W_j = 0$ provides no evidence against the null hypothesis.

⁹This is purely a problem of power and would not affect the Type I error control of knockoffs.

five SNPs, although there was substantial variance) and, by construction, upper-bounds pairwise SNP correlations by 0.5. Note that this is far from removing all dependence among the SNPs, so considering conditional independence instead of marginal dependence remains necessary for interpretation. Scientifically, we consider a selected SNP to be a true discovery if and only if it is the representative of a cluster containing a truly important SNP.

The second challenge is the one discussed in Section 3.4, and we use the approximate SDP knockoff construction proposed there. The approximate covariance matrix was just the estimated covariance matrix with zeros outside of the block diagonal, with the blocks chosen by single-linkage clustering on the estimated correlation matrix, aggregating clusters up the dendrogram as much as possible subject to a maximum cluster size of 999. In this case, even separating the problem by chromosome, the SDP construction was computationally infeasible and the equicorrelated construction produced extremely small s : $\text{mean}(s^{\text{EQ}}) = 0.08$. The parallelized approximate SDP construction took just a matter of hours to run, and increased s on average by almost an order of magnitude, with $\text{mean}(s^{\text{ASDP}}) = 0.57$.

Although it incorporates strong prior information, our estimate of the joint distribution of the SNPs is still an approximation, so our first step is to test the robustness of knockoffs to this approximation.

6.1 Simulations with genetic design matrix

Although we saw in Section 5.4 that MF knockoffs has some robustness to using an estimated covariate distribution, the situation for this data set differs from those simulations in two ways: (1) the marginal distribution of the covariates is now discrete and (2) the covariance estimator is different. One way to check the robustness of knockoffs to this particular joint distribution approximation is to take a reasonable model for $Y|X_1, \dots, X_p$ and simulate artificial response data using the real covariate data itself. If we split the rows of our design matrix into 10 smaller data sets (re-estimating the joint covariate distribution each time), for each conditional model we can run knockoffs 10 times and compute the realized FDP each time. Then averaging the 10 FDPs gives an estimate of knockoffs’ FDR for that conditional model. In an attempt to make each smaller data set have size more comparable to the $n \approx 5,000$ in our actual experiment, we combined the healthy and CD genetic information with that from 5 other diseases from the same data set:¹⁰ coronary artery disease, hypertension, rheumatoid arthritis, type 1 diabetes, and type 2 diabetes. This made for 14,708 samples altogether. To further increase the effective sample size and match the actual experiment, we used a random subset of 1000 observations for choosing cluster representatives, but the *same* 1000 for each smaller data set, so that each subsampled data set contained $\approx 1,400$ unique samples, +1000 common samples (for each of these common samples $\tilde{X}_{ij} = X_{ij}$). For computational reasons, we used only the first chromosome in this experiment, so each of our simulations had (pre-clustering) 29,258 covariates. The conditional model for the response was chosen to be a logistic regression model with 60 nonzero coefficients of random signs and locations uniformly chosen from among the original (not just representatives) SNPs. The coefficient amplitude was varied, and for each amplitude value, 10 different conditional models (different random locations of the nonzero coefficients) were simulated. Each simulation ran the exact same covariance estimation, SNP clustering and representative selection, knockoff construction, and knockoff selection procedure as used for the real data. Figure 9 shows boxplots (over conditional models) of the FDR and power at each amplitude. As hoped, the FDR is consistently controlled over a wide range of powers.

6.2 Results on real data

Encouraged by the simulation results of the previous section, we proceeded to run knockoffs with a nominal FDR level of 10% on the full genetic design matrix and real Crohn’s disease outcomes. Since knockoffs is a randomized procedure, we re-ran knockoffs 10 times (after choosing the representatives) and recorded the selected SNPs over all repetitions, summarized in Table 2. The serial computation time for a single run of knockoffs was about 6 hours, but the knockoff generation process is trivially parallelizable over chromosomes, so with 20 available computation nodes, the total parallelized computation time was about one hour. Although in this case we certainly do not know the ground truth, we can get some sort of confirmation by comparing to the results of studies with newer and much larger data sets than ours. In particular, we compared with the results of Franke et al. (2010), which used roughly 22,000 cases and 29,000 controls, or about 10 times the sample size of the WTCCC data. We also compare to the WTCCC (2007) results, where the p-value cutoff used was justified as controlling the Bayesian FDR at close to 10%—the same level we use. We consider discovered clusters in different studies to correspond (“Yes” in Table 2) if their position ranges overlap, and to nearly correspond (“Yes*” in Table 2) if the distance from our discovered cluster to the nearest cluster in the other study was less than the width of that cluster in the other study.

One thing to notice in the table is that a small number of the discovered clusters actually overlap with other clusters, specifically the clusters represented by rs11805303 and rs11209026 on chromosome 1, and rs17234657 and

¹⁰Bipolar disorder was also part of this data set, but a formatting error in the data we were given prevented us from including it.

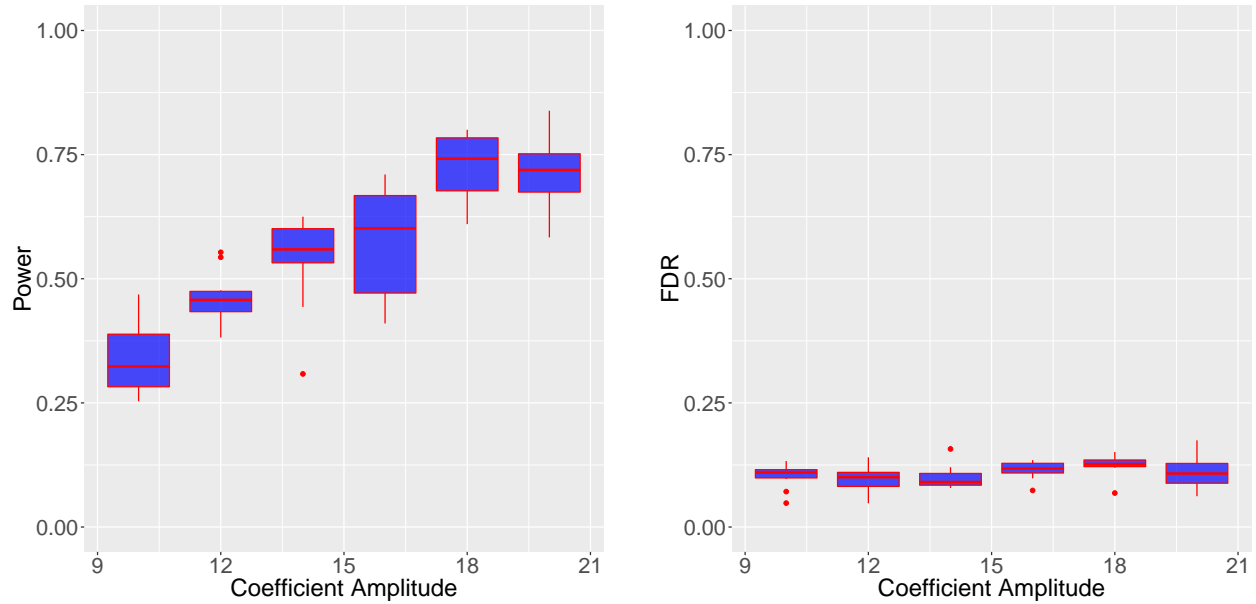


Figure 9: Power and FDR (target is 10%) for knockoffs with the LCD statistic applied to subsamples of a real genetic design matrix. Each boxplot represents 10 different logistic regression models with 60 nonzero coefficients, and for each model, 1000 common observations were used for picking cluster representatives, and the remaining 13,708 observations were divided into 10 disjoint parts and knockoffs run on each part, with power and FDR for that model then computed by averaging the results over those 10 parts.

rs7726744 on chromosome 5. And although they don't overlap one another, the three nearby clusters represented by rs3135503, rs6500315, and rs7186163 on chromosome 16 all overlap the same discovered region in Franke et al. (2010). Although puzzling at first, this phenomenon is readily explained by one of four possibilities:

- By construction, the representatives of overlapping clusters are not very highly-correlated (less than 0.5), so the fact that knockoffs chose multiple clusters in the same region may mean there are multiple important SNPs in this region, with one (or more) in each cluster. Focusing on the clusters on chromosome 1, the same region on the IL23R gene was reported in Franke et al. (2010) to have by far the strongest signal (estimated odds ratio of 2.66, next highest was 1.53) among all the regions they identified. If we conclude from this that the region or gene is fundamentally important for Crohn's disease, it stands to reason that mutations at multiple nearby but distinct loci could have important detrimental effects of their own.
- There could be an important SNP located between the two clusters, but which was not in the data. Then the two clusters on either side would both be conditionally dependent on the response, and knockoffs would be correct to reject them.
- One or more of the overlapping clusters could be mundane false discoveries caused by null covariates that happened to take values more conditionally related to the response than would be typical. This would be a facet of the data itself, and thus an unavoidable mistake.
- The tandem discoveries could also be due to a breakdown in our covariate distribution estimate. If the covariance between the two representatives were substantially underestimated and one had a large effect while the other was null, then the Lasso would have a (relatively) hard time distinguishing the two original variables, but a much easier time separating them from the knockoff of the null representative, since it is less correlated with the signal variable. As a result, the null variable and its knockoff would not be exchangeable as required by knockoffs, and a consistent error could be made. However, given that the empirical and estimated correlations between the two representatives on chromosome 1 were -0.1813 and -0.1811, respectively, and for the two on chromosome 5 were -0.2287 and -0.2286, respectively, it seems unlikely that we have made a gross estimation error. Also, note that the separation of nearly all the discovered regions, along with the simulations of the previous subsection, suggest this effect is at worst very small among our discoveries.

Overlapping clusters aside, the knockoffs results display a number of advantages over the original marginal analysis in WTCCC (2007):

- First, the power is much higher, with WTCCC (2007) making 9 discoveries, while knockoffs made 18 discoveries on average, doubling the power.
- Quite a few of the discoveries made by knockoffs that were confirmed by Franke et al. (2010) were not discovered in WTCCC (2007)’s original analysis.
- Knockoffs made a number of discoveries not found in either WTCCC (2007) or Franke et al. (2010). Of course we expect some (roughly 10%) of these to be false discoveries, particularly towards the bottom of the table. However, especially given the evidence from the simulations of the previous subsection suggesting the FDR is controlled, it is likely that many of these correspond to true discoveries. Indeed, evidence from independent studies about adjacent genes shows some of the top hits to be promising candidates. For example, the closest gene to rs6601764 is KLF6, which has been found to be associated with multiple forms of IBD, including CD and ulcerative colitis (Goodman et al., 2016); and the closest gene to rs4692386 is RBP-J, which has been linked to CD through its role in macrophage polarization (Barros et al., 2013).

Note that these benefits required relatively little customization of the knockoff procedure. For instance, WTCCC (2007) used marginal tests specifically tailored to SNP case-control data, while we simply used the LCD statistic. We conjecture that the careful use of knockoffs by domain experts would compound the advantages of knockoffs, as such users could devise more powerful statistics and better model/cluster the covariates for their particular application.

7 Discussion

This paper introduced a novel approach to variable selection in general non-parametric models, which teases apart important from irrelevant variables while guaranteeing Type I error control. This approach is a significant extension of the knockoff filter from Barber and Candès (2015) for the linear model. A distinctive feature of our approach is that selection is achieved without ever constructing p-values. This is attractive since (1) p-values are not needed and (2) it is unclear how they could be efficiently constructed, in general. (The conditional randomization approach we proposed is one way of getting such p-values but it comes at a computational cost.)

Deployment in highly correlated settings We posed a simple question: which variables does a response of interest depend upon? In many problems, there may not be enough “resolution” in the data to tell whether Y depends on X_1 or, instead, upon X_2 when the two are strongly correlated. This issue is apparent in our genetic analysis of Crohn’s disease from Section 6, where co-located SNPs may be extremely correlated. In such examples, controlling the FDR may not be a fruitful question. A more meaningful question is whether the response appears to depend upon a group of correlated variables while controlling for the effects of a number of other variables (e.g., from SNPs located in a certain region of the genome while controlling for the effects of SNPs elsewhere on the chromosomes). In such problems, we envision applying our techniques to grouped variables: one possibility is to develop a model-free group knockoff approach following Dai and Barber (2016). Another is to construct group representatives and proceed as we have done in Section 6. It is likely that there are several other ways to formulate a meaningful problem and solution.

Open questions Admittedly, this paper may pose more problems than it solves; we close our discussion with a few of them below.

- *How do we construct MF knockoffs?* Even though we presented a general strategy for constructing knockoffs, we have essentially skirted this issue except for the important case of Gaussian covariates. It would be important to address this problem, and write down concrete algorithms for some specific distributions of features of practical relevance.
- *Which MF knockoffs?* Even in the case of Gaussian covariates, the question remains of how to choose $\text{corr}(X_j, \tilde{X}_j)$ or, equivalently, the parameter s_j from Section 3 since $\text{corr}(X_j, \tilde{X}_j) = 1 - s_j$. Should we make the marginal correlations small? Should we make the partial correlations small? Should we take an information-theoretic approach and minimize a functional of the joint distribution such as the mutual information between X and \tilde{X} ?
- *What would we do with multiple MF knockoffs?* As suggested in Barber and Candès (2015), we could in principle construct multiple knockoff variables $(\tilde{X}^{(1)}, \dots, \tilde{X}^{(d)})$ in such a way that the $(d+1)p$ -dimensional family

$(X, \tilde{X}^{(1)}, \dots, \tilde{X}^{(d)})$ obeys the following extended exchangeability property: for any variable X_j , any permutation in the list $(X_j, \tilde{X}_j^{(1)}, \dots, \tilde{X}_j^{(d)})$ leaves the joint distribution invariant. On the one hand, such constructions would yield more accurate information since we could compute, among multiple knockoffs, the rank with which an original variable enters into a model. On the other hand, this would constrain the construction of knockoffs a bit more, perhaps making them less distinguishable from the original features. What is the right trade off?

Another point of view is to construct several knockoff matrices exactly as described in the paper. Each knockoff matrix would yield a selection, with each selection providing FDR control as described in this paper. Now an important question is this: is it possible to combine/aggregate all these selections leading to an increase in power while still controlling the FDR?

- *Can we prove some form of robustness?* Although our theoretical guarantees rely on the knowledge of the joint covariate distribution, our empirical studies demonstrate robustness when this distribution is simply estimated from data. For instance, the estimation of the precision matrix for certain Gaussian designs seems to have rather secondary effects on FDR and power levels. It would be interesting to provide some theoretical insights into this phenomenon.
- *Which feature importance statistics should we use?* The knockoff framework can be seen as an inference machine: the statistician provides the test statistic W_j and the machine performs inference. It is of interest to understand which statistics yield high power, as well as design new ones.
- *Can we speed up the conditional randomization testing procedure?* Conditional randomization provides a powerful alternative method for controlling the false discovery rate in model-free variable selection, but at a computational cost that is currently prohibitive for large problems. However, there exist a number of promising directions for speeding it up, including: importance sampling to estimate small p-values with fewer randomizations, faster feature statistics T_j with comparable or higher power than the absolute value of lasso-estimated coefficients, and efficient computation re-use and warm starts to take advantage of the fact that each randomization changes only a single column of the design matrix.

In conclusion, much remains to be done. On the upside, though, we have shown how to select features in high-dimensional nonlinear models (e.g., GLMs) in a reliable way. This arguably is a fundamental problem, and it is really not clear how else it could be achieved.

Acknowledgments

E. C. was partially supported by the Office of Naval Research under grant N00014-16-1-2712, and by the Math + X Award from the Simons Foundation. Y. F. was partially supported by NSF CAREER Award DMS-1150318. L. J. was partially supported by NIH training grant T32GM096982. J. L. was partially supported by a grant from the Simons Foundation. E. C. would like to thank Malgorzata Bogdan, Amir Dembo and Chiara Sabatti for helpful discussions regarding this project. E. C. would also like to thank Sabatti for superb feedback regarding an earlier version of the paper. L. J. would like to thank Kaia Mattioli for her help in understanding certain genetic principles.

References

- Affymetrix (2006). BRLMM: an improved genotype calling method for the genechip human mapping 500k array set. Technical report, Affymetrix.
- Athey, S., Imbens, G. W., and Wager, S. (2016). Efficient inference of average treatment effects in high dimensions via approximate residual balancing. *arXiv preprint arXiv:1604.07125*.
- Barber, R. F. and Candès, E. J. (2015). Controlling the false discovery rate via knockoffs. *Ann. Statist.*, 43(5):2055–2085.
- Barber, R. F. and Candès, E. J. (2016). A knockoff filter for high-dimensional selective inference. *arXiv preprint arXiv:1602.03574*.
- Barros, M. H. M., Hauck, F., Dreyer, J. H., Kempkes, B., and Niedobitek, G. (2013). Macrophage polarisation: an immunohistochemical approach for identifying m1 and m2 macrophages. *PLoS ONE*, 8(11).
- Benjamini, Y. (2010). Discovering the false discovery rate. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, 72(4):405–416.

- Benjamini, Y. and Heller, R. (2007). False discovery rates for spatial signals. *Journal of the American Statistical Association*, 102(480):1272–1281.
- Benjamini, Y. and Hochberg, Y. (1995). Controlling the false discovery rate: A practical and powerful approach to multiple testing. *Journal of the Royal Statistical Society. Series B (Methodological)*, 57(1):pp. 289–300.
- Benjamini, Y. and Yekutieli, D. (2001). The control of the false discovery rate in multiple testing under dependency. *Ann. Statist.*, 29(4):1165–1188.
- Berk, R., Brown, L., Buja, A., Zhang, K., and Zhao, L. (2013). Valid post-selection inference. *Ann. Statist.*, 41(2):802–837.
- Bogdan, M., van den Berg, E., Sabatti, C., Su, W., and Candès, E. J. (2015). SLOPE—adaptive variable selection via convex optimization. *Ann. Appl. Stat.*, 9(3):1103–1140.
- Brzyski, D., Peterson, C. B., Sobczyk, P., Candès, E. J., Bogdan, M., and Sabatti, C. (2016). Controlling the rate of gwas false discoveries. *bioRxiv*.
- Candès, E. J. and Plan, Y. (2009). Near-ideal model selection by ℓ_1 minimization. *Ann. Statist.*, 37(5A):2145–2177.
- Chouldechova, A. (2014). *False discovery rate control for spatial data*. PhD thesis, Stanford University.
- Cong, L., Ran, F. A., Cox, D., Lin, S., Barretto, R., Habib, N., Hsu, P. D., Wu, X., Jiang, W., Marraffini, L. A., and Zhang, F. (2013). Multiplex genome engineering using crispr/cas systems. *Science*, 339(6121):819–823.
- Dai, R. and Barber, R. F. (2016). The knockoff filter for fdr control in group-sparse and multitask regression. *arXiv preprint arXiv:1602.03589*.
- Doran, G., Muandet, K., Zhang, K., and Schölkopf, B. (2014). A permutation-based kernel conditional independence test. *The 30th Conference on Uncertainty in Artificial Intelligence*.
- Edwards, D. (2000). *Introduction to Graphical Modelling*. Springer New York, New York, NY.
- Efron, B. (2010). *Large-scale inference: empirical Bayes methods for estimation, testing, and prediction*, volume 1. Cambridge University Press.
- Franke, A., McGovern, D. P., Barrett, J. C., Wang, K., Radford-Smith, G. L., Ahmad, T., Lees, C. W., Balschun, T., Lee, J., Roberts, R., Anderson, C. A., Bis, J. C., Bumpstead, S., Ellinghaus, D., Festen, E. M., Georges, M., Green, T., Haritunians, T., Jostins, L., Latiano, A., Mathew, C. G., Montgomery, G. W., Prescott, N. J., Raychaudhuri, S., Rotter, J. I., Schumm, P., Sharma, Y., Simms, L. A., Taylor, K. D., Whiteman, D., Wijmenga, C., Baldassano, R. N., Barclay, M., Bayless, T. M., Brand, S., Büning, C., Cohen, A., Colombel, J.-F. F., Cottone, M., Stronati, L., Denson, T., De Vos, M., D’Inca, R., Dubinsky, M., Edwards, C., Florin, T., Franchimont, D., Gearry, R., Glas, J., Van Gossum, A., Guthery, S. L., Halfvarson, J., Verspaget, H. W., Hugot, J.-P. P., Karban, A., Laukens, D., Lawrance, I., Lemann, M., Levine, A., Libioulle, C., Louis, E., Mowat, C., Newman, W., Panés, J., Phillips, A., Proctor, D. D., Regueiro, M., Russell, R., Rutgeerts, P., Sanderson, J., Sans, M., Seibold, F., Steinhart, A. H., Stokkers, P. C., Torkvist, L., Kullak-Ublick, G., Wilson, D., Walters, T., Targan, S. R., Brant, S. R., Rioux, J. D., D’Amato, M., Weersma, R. K., Kugathasan, S., Griffiths, A. M., Mansfield, J. C., Vermeire, S., Duerr, R. H., Silverberg, M. S., Satsangi, J., Schreiber, S., Cho, J. H., Annese, V., Hakonarson, H., Daly, M. J., and Parkes, M. (2010). Genome-wide meta-analysis increases to 71 the number of confirmed Crohn’s disease susceptibility loci. *Nature genetics*, 42(12):1118–1125.
- Freedman, D. and Lane, D. (1983). A nonstochastic interpretation of reported significance levels. *Journal of Business & Economic Statistics*, 1(4):292–298.
- Frommlet, F., Ruhaltiner, F., Twaróg, P., and Bogdan, M. (2012). Modified versions of bayesian information criterion for genome-wide association studies. *Computational Statistics & Data Analysis*, 56(5):1038 – 1051.
- George, E. I. and McCulloch, R. E. (1997). Approaches for bayesian variable selection. *Statistica Sinica*, 7(2):339–373.
- Goodman, W., Omenetti, S., Date, D., Di Martino, L., De Salvo, C., Kim, G., Chowdhry, S., Bamias, G., Cominelli, F., Pizarro, T., et al. (2016). Klf6 contributes to myeloid cell plasticity in the pathogenesis of intestinal inflammation. *Mucosal immunology*.

- Guan, Y. and Stephens, M. (2011). Bayesian variable selection regression for genome-wide association studies and other large-scale problems. *Ann. Appl. Stat.*, 5(3):1780–1815.
- Haldane, J. B. S. and Waddington, C. H. (1931). Inbreeding and linkage. *Genetics*, 16(4):357–374.
- He, Q. and Lin, D.-Y. (2011). A variable selection method for genome-wide association studies. *Bioinformatics*, 27(1):1–8.
- Hoggart, C. J., Whittaker, J. C., De Iorio, M., and Balding, D. J. (2008). Simultaneous analysis of all snps in genome-wide and re-sequencing association studies. *PLoS Genet*, 4(7):1–8.
- Janson, L. and Su, W. (2016). Familywise error rate control via knockoffs. *Electron. J. Statist.*, 10(1):960–975.
- Lee, J. D., Sun, D. L., Sun, Y., and Taylor, J. E. (2016). Exact post-selection inference, with application to the lasso. *Ann. Statist.*, 44(3):907–927.
- Li, J., Das, K., Fu, G., Li, R., and Wu, R. (2011). The bayesian lasso for genome-wide association studies. *Bioinformatics*, 27(4):516–523.
- Lockhart, R., Taylor, J., Tibshirani, R. J., and Tibshirani, R. (2014). A significance test for the lasso. *Ann. Statist.*, 42(2):413–468.
- Lv, J. and Liu, J. S. (2014). Model selection principles in misspecified models. *Journal of the Royal Statistical Society Series B*, 76:141–167.
- Pacifico, M. P., Genovese, C., Verdinelli, I., and Wasserman, L. (2004). False discovery control for random fields. *Journal of the American Statistical Association*, 99(468):1002–1014.
- Pearl, J. (1988). *Probabilistic inference in intelligent systems*. Morgan Kaufmann, San Mateo, CA.
- Peters, J. M., Colavin, A., Shi, H., Czarny, T. L., Larson, M. H., Wong, S., Hawkins, J. S., Lu, C. H., Koo, B.-M., Marta, E., et al. (2016). A comprehensive, crispr-based functional analysis of essential genes in bacteria. *Cell*, 165(6):1493–1506.
- Rosenbaum, P. R. (1984). Conditional permutation tests and the propensity score in observational studies. *Journal of the American Statistical Association*, (79):pp. 565–574.
- Rosenbaum, P. R. and Rubin, D. B. (1983). The central role of the propensity score in observational studies for causal effects. *Biometrika*, 70(1):41–55.
- Saltelli, A., Ratto, M., Andres, T., Campolongo, F., Cariboni, J., Gatelli, D., Saisana, M., and Tarantola, S. (2008). *Global sensitivity analysis: the primer*. John Wiley & Sons.
- Siegmund, D. O., Zhang, N. R., and Yakir, B. (2011). False discovery rate for scanning statistics. *Biometrika*, 98(4):979–985.
- Strobl, C. and Zeileis, A. (2008). Danger: High power! ? exploring the statistical properties of a test for random forest variable importance.
- Tang, H., Coram, M., Wang, P., Zhu, X., and Risch, N. (2006). Reconstructing genetic ancestry blocks in admixed individuals. *The American Journal of Human Genetics*, 79(1):1–12.
- van de Geer, S., Bühlmann, P., Ritov, Y., and Dezeure, R. (2014). On asymptotically optimal confidence regions and tests for high-dimensional models. *Ann. Statist.*, 42(3):1166–1202.
- Wager, S. and Athey, S. (2016). Estimation and inference of heterogeneous treatment effects using random forests. *arXiv preprint arXiv:1510.04342*.
- Wen, X. and Stephens, M. (2010). Using linear predictors to impute allele frequencies from summary or pooled genotype data. *Ann. Appl. Stat.*, 4(3):1158–1182.
- WTCCC (2007). Genome-wide association study of 14,000 cases of seven common diseases and 3,000 shared controls. *Nature*, 447(7145):661–678.

Wu, T. T., Chen, Y. F., Hastie, T., Sobel, E., and Lange, K. (2009). Genome-wide association analysis by lasso penalized logistic regression. *Bioinformatics*, 25(6):714–721.

Zhao, P. and Yu, B. (2006). On model selection consistency of lasso. *Journal of Machine Learning Research*, 7(Nov):2541–2563.

A Sequential conditional independent pairs algorithm

We will prove by induction on j that Algorithm 1 produces knockoffs that satisfy the exchangeability property (3.1). We prove the result for the discrete case; the general case follows the same argument with a slightly more careful measure-theoretic treatment using Radon–Nikodym derivatives instead of probability mass functions. Below we denote the probability mass function (PMF) of $(X_{1:p}, \tilde{X}_{1:j-1})$ by $\mathcal{L}(X_{-j}, X_j, \tilde{X}_{1:j-1})$.

Induction Hypothesis. After j steps, every pair X_k, \tilde{X}_k is exchangeable in the joint distribution of $(X_{1:p}, \tilde{X}_{1:j})$ for $k = 1, \dots, j$.

By construction, the induction hypothesis is true after 1 step since X_1 and \tilde{X}_1 are conditionally independent and have the same marginal distribution (this implies conditional exchangeability). Assuming the induction hypothesis holds until $j - 1$, we prove that it holds after j steps. Note that by assumption, \mathcal{L} is symmetric in X_k, \tilde{X}_k (that is, the function value remains unchanged when the argument values X_k, \tilde{X}_k are swapped) for $k = 1, \dots, j - 1$. Also, the conditional PMF of \tilde{X}_j given $X_{1:p}, \tilde{X}_{1:j-1}$ is given by

$$\frac{\mathcal{L}(X_{-j}, \tilde{X}_j, \tilde{X}_{1:j-1})}{\sum_u \mathcal{L}(X_{-j}, u, \tilde{X}_{1:j-1})}.$$

Therefore, the joint PMF of $(X_{1:p}, \tilde{X}_{1:j})$ is given by the product of the aforementioned conditional PMF with the joint PMF of $(X_{1:p}, \tilde{X}_{1:j-1})$:

$$\frac{\mathcal{L}(X_{-j}, X_j, \tilde{X}_{1:j-1}) \mathcal{L}(X_{-j}, \tilde{X}_j, \tilde{X}_{1:j-1})}{\sum_u \mathcal{L}(X_{-j}, u, \tilde{X}_{1:j-1})}. \quad (\text{A.1})$$

Exchangeability of X_j, \tilde{X}_j follows from the symmetry of (A.1) to those two values. For $k < j$, note that (A.1) only depends on X_k, \tilde{X}_k through the function \mathcal{L} , and that \mathcal{L} is symmetric in X_k, \tilde{X}_k . Therefore, (A.1) is also symmetric in X_k, \tilde{X}_k , and therefore the pair is exchangeable in the joint distribution of $(X_{1:p}, \tilde{X}_{1:j})$.

B Bayesian knockoff statistics

The data for the simulation of Section 5.2.2 was drawn from:

$$\begin{aligned} X_j &\stackrel{\text{iid}}{\sim} \mathcal{N}(0, 1/n), \quad j \in \{1, \dots, p\}, \\ \beta_j &\stackrel{\text{iid}}{\sim} \mathcal{N}(0, \tau^2), \quad j \in \{1, \dots, p\}, \\ \delta_j &\stackrel{\text{iid}}{\sim} \text{Bernoulli}(\pi), \quad j \in \{1, \dots, p\}, \\ \frac{1}{\sigma^2} &\sim \text{Gamma}(A, B) \quad (\text{shape/scale parameterization, as opposed to shape/rate}), \\ Y &\sim \mathcal{N}\left(\sum_{j:\delta_j=1} X_j \beta_j, \sigma^2\right). \end{aligned}$$

The simulation used $n = 300$, $p = 1000$, with parameter values $\pi = \frac{60}{1000}$, $A = 5$, $B = 4$, and τ varied along the x-axis of the plot.

To compute the Bayesian variable selection (BVS) knockoff statistic, we used a Gibbs sampler on the following

model (treating X_1, \dots, X_p and $\tilde{X}_1, \dots, \tilde{X}_p$ as fixed):

$$\begin{aligned}
\beta_j &\stackrel{\text{iid}}{\sim} \mathcal{N}(0, \tau^2), \quad j \in \{1, \dots, p\}, \\
\tilde{\beta}_j &\stackrel{\text{iid}}{\sim} \mathcal{N}(0, \tau^2), \quad j \in \{1, \dots, p\}, \\
\lambda_j &\stackrel{\text{iid}}{\sim} \text{Bernoulli}(\pi), \quad j \in \{1, \dots, p\}, \\
(\delta_j, \tilde{\delta}_j) &\stackrel{\text{iid}}{\sim} \left\{ \begin{array}{ll} (0, 0) & \text{if } \lambda_j = 0 \\ (0, 1) \text{ w.p. } 1/2 & \text{if } \lambda_j = 1 \\ (1, 0) \text{ w.p. } 1/2 & \text{if } \lambda_j = 1 \end{array} \right\}, \quad j \in \{1, \dots, p\}, \\
\frac{1}{\sigma^2} &\sim \text{Gamma}(A, B) \quad (\text{shape/scale parameterization, as opposed to shape/rate}), \\
Y &\sim \mathcal{N} \left(\sum_{j:\delta_j=1} X_j \beta_j + \sum_{j:\tilde{\delta}_j=1} \tilde{X}_j \tilde{\beta}_j, \sigma^2 \right),
\end{aligned}$$

which requires only a very slight modification of the procedure in George and McCulloch (1997). After computing the posterior probabilities $\hat{\delta}_j$ and $\hat{\tilde{\delta}}_j$ with 500 Gibbs samples (after 50 burn-in samples), we computed the j th knockoff statistic as

$$W_j = \hat{\delta}_j - \hat{\tilde{\delta}}_j.$$

C WTCCC data

The SNP arrays came from an Affymetrix 500K chip, with calls made by the BRLMM algorithm Affymetrix (2006). SNPs satisfying any of the following conditions were removed:

- minor allele frequency (MAF) $< 1\%$,
- Hardy–Weinberg equilibrium test p-value $< 0.01\%$,
- missing $> 5\%$ of values (values were considered missing if the BRLMM score was > 0.5),
- position was listed as same as another SNP (this occurred just for the pair rs16969329 and rs4886982; the former had smaller MAF and was removed),
- position was not in genetic map.

Furthermore, subjects missing $> 5\%$ of values were removed. Following the original paper describing/using this data (WTCCC, 2007), we did not adjust for population structure.

Any missing values that remained after the above preprocessing were replaced by the mean of the nonmissing values for their respective SNPs.

D Knockoffs with selection

First, recall that the results of Theorem 3.4 hold if for any subset $S \subset \mathcal{H}_0$, we have

$$([\mathbf{X}, \tilde{\mathbf{X}}]_{\text{swap}(S)}, y) \stackrel{d}{=} ([\mathbf{X}, \tilde{\mathbf{X}}], y). \quad (\text{D.1})$$

In fact, MF knockoffs are defined in such a way that this property holds. Now, the procedure employed in Section 6 to construct knockoffs is slightly different from that described in the rest of the paper. Explicitly, the data looks like

$$\mathbf{X} = \begin{bmatrix} \mathbf{X}^{(1)} \\ \mathbf{X}^{(2)} \end{bmatrix}, \quad y = \begin{bmatrix} y^{(1)} \\ y^{(2)} \end{bmatrix},$$

where $y^{(1)}$ is 983×1 , $\mathbf{X}^{(1)}$ is $983 \times 71,145$, $y^{(2)}$ is 3930×1 and $\mathbf{X}^{(2)}$ is $3930 \times 71,145$; recall that the set of samples labeled (1) is used to select the cluster representatives, and that the two sets (1) and (2) are independent of each

other. The knockoffs $\tilde{\mathbf{X}}^{(2)}$ for $\mathbf{X}^{(2)}$ are generated as described in the paper (using the ASDP construction) and for (1), we set $\tilde{\mathbf{X}}^{(1)} = \mathbf{X}^{(1)}$. To verify (D.1), it suffices to show that for any subset $S \in \mathcal{H}_0$ and each $i \in \{1, 2\}$,

$$([\mathbf{X}^{(i)}, \tilde{\mathbf{X}}^{(i)}], y^{(i)}) \stackrel{d}{=} ([\mathbf{X}^{(i)}, \tilde{\mathbf{X}}^{(i)}]_{\text{swap}(S)}, y^{(i)}) \quad (\text{D.2})$$

since the sets (1) and (2) are independent. (D.2) holds for $i = 2$ because we are following the classical knockoffs construction. For $i = 1$, (D.2) actually holds with exact equality, and not just equality in distribution.

We can also argue from the perspective of MF knockoffs from Definition 3.1. By construction, it is clear that $\tilde{\mathbf{X}}^{(2)}$ are valid MF knockoffs for $\mathbf{X}^{(2)}$. We thus study $\tilde{\mathbf{X}}^{(1)}$: since $\tilde{\mathbf{X}}^{(1)} = \mathbf{X}^{(1)}$, the exchangeability property is trivial; also, $\tilde{\mathbf{X}}^{(1)} \perp\!\!\!\perp y^{(1)} \mid \mathbf{X}^{(1)}$ since $\mathbf{X}^{(1)}$ determines $\tilde{\mathbf{X}}^{(1)}$.

Selection frequency	Cluster Representative (Cluster Size)	Chrom.	Position Range (Mb)	Confirmed in Franke et al. (2010)?	Selected in WTCCC (2007)?
100%	rs11805303 (16)	1	67.31–67.46	Yes	Yes
100%	rs11209026 (2)	1	67.31–67.42	Yes	Yes
100%	rs6431654 (20)	2	233.94–234.11	Yes	Yes
100%	rs6601764 (1)	10	3.85–3.85	No	No
100%	rs7095491 (18)	10	101.26–101.32	Yes	Yes
90%	rs6688532 (33)	1	169.4–169.65	Yes	No
90%	rs17234657 (1)	5	40.44–40.44	Yes	Yes
90%	rs3135503 (16)	16	49.28–49.36	Yes	Yes
80%	rs9783122 (234)	10	106.43–107.61	No	No
80%	rs11627513 (7)	14	96.61–96.63	No	No
60%	rs4437159 (4)	3	84.8–84.81	No	No
60%	rs7768538 (1145)	6	25.19–32.91	Yes	No
60%	rs6500315 (4)	16	49.03–49.07	Yes	Yes
60%	rs2738758 (5)	20	61.71–61.82	Yes	No
50%	rs7726744 (46)	5	40.35–40.71	Yes	Yes
50%	rs4246045 (46)	5	150.07–150.41	Yes	Yes
50%	rs2390248 (13)	7	19.8–19.89	No	No
50%	rs7186163 (6)	16	49.2–49.25	Yes	Yes
40%	rs10916631 (14)	1	220.87–221.08	No	No
40%	rs4692386 (1)	4	25.81–25.81	No	No
40%	rs7655059 (5)	4	89.5–89.53	No	No
40%	rs7759649 (2)	6	21.57–21.58	Yes*	No
40%	rs1345022 (44)	9	21.67–21.92	No	No
30%	rs6825958 (3)	4	55.73–55.77	No	No
30%	rs9469615 (2)	6	33.91–33.92	Yes*	No
30%	rs4263839 (23)	9	114.58–114.78	Yes	No
30%	rs2836753 (5)	21	39.21–39.23	No	No
10%	rs459160 (2)	1	44.75–44.75	No	No
10%	rs6743984 (23)	2	230.91–231.05	Yes	No
10%	rs2279980 (20)	5	57.95–58.07	No	No
10%	rs4959830 (11)	6	3.36–3.41	Yes	No
10%	rs13230911 (9)	7	1.9–2.06	No	No
10%	rs7807268 (5)	7	147.65–147.7	No	No
10%	rs2147240 (1)	9	71.83–71.83	No	No
10%	rs10761659 (53)	10	64.06–64.41	Yes	Yes
10%	rs4984405 (3)	15	93.06–93.08	No	No
10%	rs17694108 (1)	19	38.42–38.42	Yes	No
10%	rs3932489 (30)	20	15.01–15.09	No	No

Table 2: SNP clusters discovered to be important for Crohn’s disease over 10 repetitions of knockoffs. Clusters not found in Franke et al. (2010) represent promising sites for further investigation, especially rs6601764 and rs4692386, whose nearest genes have been independently linked to CD. See text for detailed description. SNP positions are as listed in the original data, which uses Human Genome Build 35.